# DEEPWEIGHTFLOW: RE-BASINED FLOW MATCHING FOR GENERATING NEURAL NETWORK WEIGHTS

**Anonymous authors**Paper under double-blind review

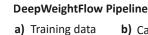
#### **ABSTRACT**

Building efficient and effective generative models for neural network weights has been a research focus of significant interest that faces challenges posed by the high-dimensional weight spaces of modern neural networks and their symmetries. Several prior generative models are limited to generating partial neural network weights, particularly for larger models, such as ResNet and ViT. Those that do generate complete weights struggle with generation speed or require finetuning of the generated models. In this work, we present DeepWeightFlow, a Flow Matching model that operates directly in weight space to generate diverse and highaccuracy neural network weights for a variety of architectures, neural network sizes, and data modalities. The neural networks generated by DeepWeightFlow do not require fine-tuning to perform well and can scale to very large networks. We apply Git Re-Basin and TransFusion for neural network canonicalization in the context of generative weight models to account for the impact of neural network permutation symmetries and to improve generation efficiency for larger model sizes. The generated networks excel at transfer learning, and ensembles of hundreds of neural networks can be generated in minutes, far exceeding the efficiency of diffusion-based methods. DeepWeightFlow models pave the way for more efficient generation of diverse sets of large neural networks.

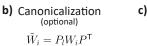
### 1 Introduction

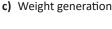
Generating neural network weights is a sampling challenge that explores the underlying high-dimensional distribution of weights, where neural networks trained on similar datasets and tasks exhibit statistical regularities. The development of generative models capable of learning the distributional properties of trained weights faces challenges of symmetries and high-dimensionality of the weight spaces. Treating large collections of neural network weights as a structured and high-dimensional data modality promises advances in model editing (Mitchell et al., 2022; Meng et al., 2022), accelerating transfer learning (Knyazev et al., 2021; Schürholt et al., 2022), facilitating uncertainty quantification (Lakshminarayanan et al., 2017), and advancing neural architecture search (Chen et al., 2019; Chen, 2023). Unlike traditional machine learning tasks that aim to optimize weights for specific downstream tasks, this concept advocates sampling from the weight space itself. In this work, we focus on the efficient generation of complete neural network weights that can achieve high performance for a given task and excel at transfer learning thus addressing fundamental limitations in current deep learning workflows, such as computational bottlenecks in iterative training, vulnerability to adversarial attacks (Goodfellow et al., 2015; Madry et al., 2018) and privacy concerns arising from training data reconstructions (Nasr et al., 2019; Tramer et al., 2022).

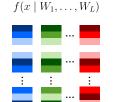
Generating neural network weights faces three main challenges: *Firstly*, neural network weights have a rich class of symmetries (Hecht-Nielsen, 1990; Entezari et al., 2022; Navon et al., 2023; Zhao et al., 2025), i.e., transformations of the weights that leave the neural network functionally invariant. Most prominently, joint permutations of hidden neurons in adjacent layers of multi-layer perceptrons (MLP) do not change the encoded function. Other architectural choices, such as incorporating attention heads or the choice of non-linear activation, can induce additional symmetries. Techniques for dealing with weight space symmetries fall into three main categories: (1) data augmentation, (2) equivariant architectures, and (3) canonicalization. Prior work, such as Wortsman et al. (2022); Wang et al. (2024); Soro et al. (2025); Saragih et al. (2025a), does not actively account for symmetries in their generative models, while others, such as Saragih et al. (2025b), use equivariant architec-



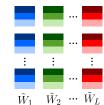
Target networks







 $W_L$ 



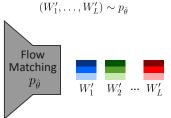


Figure 1: Schematic depiction of DeepWeightFlow. **a**) We construct a training dataset of weights by fully training neural networks with weights  $W_1, \ldots, W_L$  on a given target task. **b**) Optionally, we use canonicalization, i.e., choosing a canonical representative  $\tilde{W}_i$  from the same orbit as  $W_i$ , to break the permutation symmetry in parameter space. **c**) We train a flow model  $p_{\hat{\theta}}$  for efficient generation of high-performance weights  $(W_1, \ldots, W_L) \sim p_{\hat{\theta}}$  for the target task.

tures. Data augmentation has also been explored in weight representation learning (Schürholt et al., 2024; Shamsian et al., 2023; 2024), and to a lesser extent in weight generation (Schürholt et al., 2024; Wang et al., 2025). Finally, canonicalization has recently found application in weight space learning (Schürholt et al., 2024; Wang et al., 2024; 2025), borrowing ideas from model merging and alignment (Ainsworth et al., 2023; Rinaldi et al., 2025).

Secondly, neural network weights are high-dimensional, varying from tens of millions for a small ResNet (He et al., 2016) to hundreds of billions for modern large language models (Touvron et al., 2023; Guo et al., 2025). This challenge is often addressed by non-linear, dimensionality reduction techniques, including variational autoencoders (VAEs) (Soro et al., 2025) and graph autoencoders (Schürholt et al., 2022; Saragih et al., 2025b; Soro et al., 2025). Despite increasing efficiency, dimensionality reduction requires training an additional model for dimensionality reduction and can be detrimental to the quality of the generated weights if the compression is lossy.

*Lastly*, generative models proposed recently either generate partial weights for large models, or require finetuning post-generation, or have long generation time per sample, making them impractical.

To address these challenges, we propose DeepWeightFlow, a method for efficient generation of high-performance neural network weights via Flow Matching (FM) and apply it to MLP for vision and tabular data, as well as ResNet, and ViT (Dosovitskiy et al., 2021). We rely on canonicalization techniques, such as Git Re-Basin (Ainsworth et al., 2023) and TransFusion (Rinaldi et al., 2025), to resolve parameter permutation symmetries, and show that canonicalization aids weight generation for large neural networks but offers limited benefits when the weight space dimension is moderate. We show that neural networks generated by DeepWeightFlow excel at the target task and are competitive with state-of-the-art weight generation methods such as RPG (Wang et al., 2025), D2NWG (Soro et al., 2025), FLoWN (Saragih et al., 2025b), and P-diff (Wang et al., 2024) while overcoming several of the limitations of these models. A schematic of our methods is shown in Figure 1. While DeepWeightFlow samples directly from weight spaces, we show that the models can easily scale to generating larger networks using PCA while keeping the training and the generation time low. *In summary, the contributions of this work are as follows:* 

- DeepWeightFlow is a new method for *complete* neural network weight generation based on FM, unconditioned by dataset characteristics, task descriptions, or architectural specifications.
   DeepWeightFlow does not require additional training of autoencoders for dimensionality reduction and can scale to high-dimensional weight spaces using PCA.
- We show that our method can generate weights for neural networks with up to  $\mathcal{O}(10M)$  parameters, and diverse architectures, such as MLP, ResNet, and ViT, that, without fine-tuning, exhibit high performance on tasks from the vision and tabular domains.
- We empirically elucidate the role of parameter symmetry for weight generation, showing that
  canonicalization of the training data aids the generation of very high-dimensional weights but
  offers no additional benefit for weights of modest dimension.
- DeepWeightFlow with a simple MLP implementation, without any equivariant architecture, is far more efficient in generating diverse samples compared to diffusion-based models.

## 2 RELATED WORK

**HyperNetworks:** Early explorations of neural network generation focus on HyperNetworks, which learn neural network parameters as a relaxed temporal weight sharing process (Ha et al., 2017). HyperNetworks have been applied to generating weights through density sampling, GAN, and diffusion methods by learning latent representations of neural network weights (Ha et al., 2017; Frankle & Carbin, 2019; Ratzlaff & Fuxin, 2019; Schürholt et al., 2022; Kiani et al., 2024). They have also been used to build meta-learners – augmentations or substitutes for Stochastic Gradient Descent optimization, which condition generation of new weight checkpoints on prior weights and task losses (Peebles et al., 2022; Zhang et al., 2024a; Wang et al., 2025).

Generative Models for Neural Network Weights: Diffusion-based generative models for weights have been successful at neural network weight generation, but often do not directly resolve weight space symmetries. These approaches either provide no treatment (Wang et al., 2024), or rely on Variational Auto Encoding (VAE) methods to concurrently resolve weight symmetries and reduce the dimensionality of the generative task (Ha et al., 2017; Frankle & Carbin, 2019; Schürholt et al., 2022; Kiani et al., 2024; Soro et al., 2025). In contrast, weight canonicalization is done as a pretraining step in SANE (Schürholt et al., 2024), which uses kernel density sampling of hypernetwork latents to autoregressively populate models layer-wise, allowing for *complete* weight generation, but requires fine-tuning, unlike DeepWeightFlow. Diffusion has been applied directly to generating partial (Wang et al., 2024) or complete weights (Soro et al., 2025; Wang et al., 2025). RPG (Wang et al., 2025) generates complete weights by using a recurrent diffusion model. However, RPG shows long generation times, often taking hours to generate a set of networks that DeepWeightFlow takes minutes to complete. Subsequent Conditional Flow Matching (CFM) methods (Saragih et al., 2025b;a) explore dataset embeddings as conditioning for transfer learning and weight generation. These CFMs also report using VAE methods to reduce the dimensionality of the generative task and to resolve weight symmetries (Saragih et al., 2025b;a). We develop this further with DeepWeightFlow, which operates directly in deep weight space to generate complete weight sets, and demonstrates the viability of PCA as a strategy for surpassing  $\mathcal{O}(10M)$  parameter sets.

**Permutation Symmetries in Weight Space:** SANE (Schürholt et al., 2024) applies Git Re-Basin as a canonicalization for hypernetwork training (Schürholt et al., 2022; 2024; Ainsworth et al., 2023). Unlike DeepWeightFlow, SANE tokenizes weights layer-wise and autoregressively samples them to populate new neural models. RPG (Wang et al., 2025) uses a different strategy to address permutation symmetry by one-hot encoding models to differentiate between potential permutations of similar weights. D2NWG (Soro et al., 2025) and FLoWN (Saragih et al., 2025b) evaluate VAEs, and permutation invariant graph autoencoding and pruning methods to appeal to the manifold and lottery ticket hypotheses (Ha et al., 2017; Frankle & Carbin, 2019; Schürholt et al., 2022; Kiani et al., 2024). DeepWeightFlow extends the canonicalization methods from previous works to transformers through TransFusion, and thoroughly evaluates the impact of canonicalization on generating *complete* weight sets (Schürholt et al., 2024; Wang et al., 2024; 2025; Soro et al., 2025).

## 3 BACKGROUND

DeepWeightFlow is an FM model using an MLP architecture trained on canonicalized neural networks. In this section, we give a brief overview of the various methods we use to build it.

# 3.1 FLOW MATCHING

Flow Matching (Lipman et al., 2023) is a generative technique for learning a vector field to transport a noise vector to a target distribution. Given an unknown data distribution q(x), we define a probability path  $p_t$  for  $t \in [0,1]$  with  $p_0 \sim \mathcal{N}(0,1)$  and  $p_1 \approx q(x)$ . FM learns a vector field with parameters  $\theta$ ,  $v_{\theta}(x,t)$ , that transports  $p_0$  to  $p_1$  by minimizing

$$\mathcal{L}_{FM}(\theta) = \mathbb{E}_{t \sim \mathcal{U}[0,1], x \sim p_t(x)} \left[ \|v_{\theta}(x,t) - u(x,t)\|^2 \right], \tag{1}$$

where u(x,t) is the true vector field generating  $p_t(x)$ , and  $\mathcal{U}[0,1]$  denotes the uniform distribution on the unit interval [0,1]. This loss is minimized if  $v_\theta$  matches u, effectively following the probability path from  $p_0$  to  $p_1$ . FM offers several advantages over diffusion for neural network weight generation as it enables simpler and faster sampling, relies on direct vector field regression for training,

and scales efficiently to high-dimensional spaces, making it particularly well-suited for generating complete neural network weights.

#### 3.2 PERMUTATION SYMMETRIES OF NEURAL NETWORKS AND RE-BASIN

Permutation symmetry is a common weight space symmetry in neural networks (Hecht-Nielsen, 1990). Consider the activations  $z_\ell \in \mathbb{R}^{d_\ell}$  at the  $\ell^{\text{th}}$  layer of a simple MLP, with weights  $W_\ell \in \mathbb{R}^{d_{\ell+1} \times d_\ell}$ , biases  $b_\ell \in \mathbb{R}^{d_{\ell+1}}$ , and activation  $\sigma$ ,  $z_{\ell+1} = \sigma(W_\ell z_\ell + b_\ell)$ , where  $z_0 = x$  is the input data. Applying a permutation matrix  $P \in \mathbb{R}^{d_{\ell+1} \times d_{\ell+1}}$  of appropriate dimension, yields

$$z_{\ell+1} = P^{\mathsf{T}} P z_{\ell+1} = P^{\mathsf{T}} P \sigma(W_{\ell} z_{\ell} + b_{\ell}) = P^{\mathsf{T}} \sigma(P W_{\ell} z_{\ell} + P b_{\ell}), \tag{2}$$

where  $P^{\mathsf{T}}P = I$ . This shows that a permutation of the output features of the  $\ell^{th}$  layer, when met with the appropriate permutation of the input features of the next layer  $\ell + 1$ , will leave the overall MLP functionally invariant (Ainsworth et al., 2023).

Similar permutation symmetries (Lim et al., 2024) exist for the channels of convolutional neural networks and the attention heads of the transformer architecture (Hecht-Nielsen, 1990; Ainsworth et al., 2023; Rinaldi et al., 2025). These symmetries shape the loss landscape (Pittorino et al., 2022), impacting optimization(Neyshabur et al., 2015a; Liu, 2023; Zhao et al., 2024), generalization(Neyshabur et al., 2015b; Dinh et al., 2017), and model complexity (Zhao et al., 2025). They also impact the ability of generative models to learn distributions over neural network weights. Permutation symmetry gives rise to a highly multi-modal loss surface leading to a large set of weights that render the resulting models equivalent in task performance (Hecht-Nielsen, 1990; Lim et al., 2024).

In model alignment, weights are aligned with respect to a reference model to produce unique 'canonical' representations for each equivalence class of the weight permutation symmetry. The Git Re-Basin (Ainsworth et al., 2023) weight matching approach permutes the hidden units of an MLP such that the inner product between reference and permuted weights is maximized. The resulting optimization problem is a sum of bilinear assignment problems (SOBLAP). Git Re-Basin solves this problem approximately, using coordinate descent, reducing each layer's subproblem to a linear assignment and iterating until convergence. TransFusion (Rinaldi et al., 2025) extends this idea of weight alignment to transformers where permutation symmetries exist both in MLPs and within and between attention heads, applying iterative alignment steps to reconcile permutations of heads and hidden units. More details on this can be found in Appendix A and Appendix B.

### 4 Methods

We implement a simple MLP-based FM model. The explicit encoding of the symmetries of the neural networks is done using TransFusion for transformers and Git Re-Basin for all other architectures. We generate neural network weights independently trained from random initialization and not drawn from a sequence of checkpoints from training a single neural network, thus increasing the diversity of the training set, for training all DeepWeightFlow models.

Flow Matching Architecture and Training: DeepWeightFlow is a flow-based model Lipman et al. (2023). We use a time-conditioned neural network that predicts a velocity vector along a trajectory between source and target network weights. The source is a distribution of Gaussian noise given by  $x_0 \sim \mathcal{N}(0, \sigma^2 I)$ , and the target is a distribution of trained weights  $(x_1 \sim p_{\text{target}})$ . The source distribution has the same dimensions as the target. Given a sampled time  $t \in [0, 1]$  (uniformly distributed), an interpolated point along the straight-line trajectory is computed as  $\mu_t = (1-t)x_0 + tx_1$ . To stabilize training, stochastic points are generated by adding Gaussian noise  $x_t = \mu_t + \epsilon$ , with  $\epsilon \sim \mathcal{N}(0, \sigma^2 I)$ . The instantaneous target velocity along this linear trajectory is  $u_t = x_1 - x_0$  (since  $\frac{d\mu_t}{dt} = x_1 - x_0$ ), which is constant along the straight-line path. The network sees  $x_t$  as input, while  $u_t$  is derived from the endpoints  $(x_0, x_1)$ . The scalar time t is embedded into a higher-dimensional vector  $t_{\text{embed}} = \text{MLP}(t) \in \mathbb{R}^{d_{\text{time}}}$ , where  $d_{\text{time}}$  varies depending on the complexity of the model for which we are training DeepWeightFlow. We use a shallow MLP with layer normalization, dropout regularization, and GELU activations. This  $t_{\text{embed}}$  is concatenated with  $x_t$  and fed into the main network, allowing the network to condition on time in a learnable, flexible manner. The network maps  $(x_t, t_{\text{embed}}) \mapsto v_\theta(x_t, t)$ , where  $v_\theta$  is the learned vector field. The main network

consists of fully connected layers with LayerNorm, GELU activations, and Dropout, ending with a linear layer mapping back to the flattened weight dimension. Finally, new weight configurations are generated by integrating the learned vector field from random Gaussian inputs in the same flattened weight space as the source distribution. This integration is performed using a fourth-order Runge-Kutta (RK4) method, which ensures high-accuracy trajectories in weight space. Concretely, at each integration step, the vector field is evaluated at the current point and time, and RK4 increments are computed to update the weights. This procedure allows sampling of realistic neural network weight configurations that smoothly interpolate between source and target distributions.

Canonicalization: We apply canonicalization to align the training set to a single reference, as neural network loss landscapes are inherently degenerate due to permutation symmetries in the weight space. This simplifies the learning process without the need for complex equivariant architectures. To implement canonicalization for smaller MLPs and ResNets, we use the weight-matching procedure of Git Re-Basin (Ainsworth et al., 2023) for 100 iterations. For ViTs, we use the TransFusion procedure (Rinaldi et al., 2025) for 10 iterations as the latter uses spectral decomposition and is slower than Git Re-Basin. The detailed description of these methods can be found in Appendix A and Appendix B.

Batch Normalization Statistics Based Recalibration: We implement a post-generation recalibration procedure where batch normalization (BN) (Ioffe & Szegedy, 2015) statistics are recomputed using the training dataset for each set of generated weights. Neural networks with BN pose challenges for weight generation, as even perfectly generated weights can underperform if BN statistics are misaligned. DeepWeightFlow addresses this by recalibrating BN statistics after weight generation, ensuring models are accurate. While the FM framework successfully learns BN weight parameters ( $\gamma$  and  $\beta$ ), the running statistics (mean and variance) require more careful processing. These statistics are intrinsically tied to the training data distribution and must be precisely calibrated for each generated weight set. Our experiments, summarized in Table 6, reveal that directly transferring running statistics from a reference model yields suboptimal performance. We provide our recalibration algorithm in Algorithm 1 (Wortsman et al., 2021; 2022). Layer normalization (Ba et al., 2016) is permutation invariant and does not need recalibration (Ainsworth et al., 2023).

Training Data Generation: All training data used in this work was generated *ab initio* from a set of randomly initialized neural networks trained separately, thus generating a diverse set of neural networks. Details of the training dataset generation can be found in Appendix G. We test Deep-WeightFlow on diverse tasks such as the Iris (Fisher, 1936), MNIST (Lecun et al., 1998), Fashion-MNIST (Xiao et al., 2017), and CIFAR-10 (Krizhevsky et al., 2009) datasets. Recent work by Zeng et al. (2025) has raised concerns about the lack of diversity of weights sampled from generative models trained on checkpoints from training a single neural network (Wang et al., 2024). As described above, we do not use checkpoints from a single training run. We provide code to generate the training dataset in https://github.com/anonymousacademicc/DeepWeightFlow-ICLR and hyperparameters in Table 9 and Table 10 of Appendix G.

## 5 EXPERIMENTS

We conduct a series of experiments to evaluate the effectiveness of our approach across different architectures, training conditions, and downstream tasks. We show that DeepWeightFlow generates complete weights for MLPs, ResNets, and Vits with high accuracy, and canonicalization improves performance at low flow matching model capacity. We see that incremental PCA enables scaling DeepWeightFlow to tens of millions of parameters without accuracy degradation. Our approach is robust across diverse initialization schemes, including Kaiming, Xavier, Gaussian, and Uniform. We see that Gaussian source distributions outperform Kaiming, with variance choice being most critical at low capacity. Generated CIFAR-10 models transfer effectively to STL-10 and SVHN, both in zero-shot and fine-tuning settings. Lastly, the generated models are diverse while maintaining strong accuracy, and training and sampling are significantly faster than diffusion models such as RPG.

5.1 COMPLETE WEIGHT GENERATION ACROSS ARCHITECTURES

DeepWeightFlow generates complete neural network weights and the generated networks perform as well as the training set. In Table 1, Table 2, and Table 3, we highlight the results of generating MLPs, ResNet-18/20s and ViTs from DeepWeightFlow models. We have conducted our experiments on MNIST, Fashion-MNIST, CIFAR-10, STL-10 (Coates et al., 2011), and SVNH (Goodfel-

Table 1: Comparison of DeepWeightFlow with other SOTA neural network weight generating methods for complete generation of weights for MNIST classifiers, without finetuning.

Model	Neural Network	Original	Generated	Reference
DeepWeightFlow (w/ Git Re-Basin) DeepWeightFlow (w/o Git Re-Basin)	3-Layer MLP	$96.32 \pm 0.20$	$96.17 \pm 0.31$ $96.19 \pm 0.27$	
WeightFlow (Geometric, aligned + OT)	3-Layer MLP	93.3	78.6	Erdogan (2025)
FLoWN (Unconditioned)	medium-CNN	92.76	83.58	Saragih et al. (2025b)

Table 2: Comparison of DeepWeightFlow with other SOTA neural network weight generating models for complete ResNet-18 CIFAR-10 classifier weight generation, without fine tuning.

Model	Original	Generated (Partial)	Generated (Complete)	Reference Reference
DeepWeightFlow (w/ Git Re-Basin) DeepWeightFlow (w/o Git Re-Basin)	$94.45 \pm 0.14$		$93.55 \pm 0.13$ $93.47 \pm 0.20$	
RPG	95.3	-	95.1	Wang et al. (2025)
SANE	$92.14 \pm 0.12$	-	$68.6 \pm 1.2$	Schürholt et al. (2024)
D2NWG	94.56	-	$94.57 \pm 0.0$	Soro et al. (2025)
P-diff (best neural network)	94.54	94.36	-	Wang et al. (2024) (Saragih et al., 2025b)
FLoWN (best neural network)	94.54	94.36	_	Saragih et al. (2025b)

Table 3: Comparison of DeepWeightFlow with other SOTA neural network weight generating models for ViT family CIFAR-10 classifiers, without finetuning. We have used ViT-small-192, indicating an embedding dimension of 192 Wang et al. (2025); Schürholt et al. (2024); Soro et al. (2025); Dosovitskiy et al. (2021).

Model	neural network	Original	Generated	Reference
DeepWeightFlow (w/ TransFusion) DeepWeightFlow (w/o TransFusion)	Vit-Small-192	$83.30 \pm 0.29$	$83.07 \pm 0.42$ $82.58 \pm 0.07$	
P-diff (Best)	ViT-mini	73.0	73.6	Wang et al. (2024)
RPG	ViT-Base	98.7	98.9	Wang et al. (2025)

low et al., 2013) datasets. As noted before, we generate the *complete* weights for all neural networks, including those with batch normalization such as ResNet-18 and ResNet-20. The comprehensive weight generation scope of DeepWeightFlow is unlike existing approaches such as FLoWN (Saragih et al., 2025b) and P-diff (Wang et al., 2024), which primarily generate only partial weight sets (limited to batch normalization parameters due to lack of scalability with neural network parameter size). Moreover, DeepWeightFlow generated networks perform as well as the training set without the requirement of additional conditioning during training or inference. With sufficient flow model capacity, performance converges regardless of canonicalization or noise scheduling strategy, suggesting that model capacity can compensate for suboptimal design choices. The choice of source distribution significantly impacts FM performance and generated model diversity (cf. Figure 2).

Effect of Source Distributions: Critical to the success of DeepWeightFlow, is the careful selection of the standard deviation parameter of the source distribution: optimal results are achieved when the source distribution's standard deviation matches or slightly undershoots that of the target weight distribution. Our empirical analysis demonstrates that Gaussian noise consistently outperforms alternative initializations (e.g., Kaiming initialization) as the source distribution (Table 8 in Appendix F). This sensitivity is particularly pronounced in smaller flow models, where insufficient capacity amplifies the importance of proper initialization (Saragih et al., 2025b).

**Scaling with PCA:** DeepWeightFlow can scale to large neural networks using PCA. For models with tens of millions of parameters, we employ incremental PCA to reduce the dimensionality of flattened weight vectors before FM, followed by inverse transformation post-generation. This approach maintains accuracy levels, as can be seen from Table 7 in Appendix E, while enabling tractable training of DeepWeightFlow for large-scale architectures. This demonstrates the feasibil-

ity of extending our methodology to generate complete weight sets for contemporary large neural networks without the requirement of training additional models for dimensionality reduction, such as autoencoders, as is often done for latent diffusion-based models (Wang et al., 2024).

Table 4: Canonicalization is beneficial when DeepWeightFlow has limited capacity, leading to superior performance. As model capacity increases, both canonicalized and non-canonicalized models perform comparably, with the best results highlighted in bold.

			Original	Gene	erated
Dataset	Architecture	${d_h}^*$	(accuracy) mean $\pm$ st. dev.	with re-basin (accuracy) mean $\pm$ st. dev.	without re-basin (accuracy) mean $\pm$ st. dev.
Iris	MLP	256 128 64 32	$90.70 \pm 2.02$	$91.43 \pm 2.07$ $91.43 \pm 2.46$ $91.87 \pm 2.23$ $90.80 \pm 2.54$	$91.03 \pm 2.20$ $90.87 \pm 3.25$ $90.80 \pm 4.86$ $88.93 \pm 6.09$
MNIST	MLP	512 256 128 64	$96.32 \pm 0.20$	$96.17 \pm 0.31$ $96.21 \pm 0.28$ $91.74 \pm 10.37$ $57.80 \pm 9.85$	$96.19 \pm 0.27$ $96.20 \pm 0.23$ $89.71 \pm 17.93$ $25.54 \pm 12.90$
Fashion-MNIST	MLP	512 256 128 64	$89.24 \pm 0.27$	$89.10 \pm 0.29$ $89.06 \pm 0.29$ $88.09 \pm 2.24$ $77.76 \pm 3.72$	89.11 ± 0.28 89.02 ± 0.30 85.81 ± 11.32 53.35 ± 30.49
CIFAR-10	ResNet-20	512 256 128 64	$73.62 \pm 2.24$	$75.07 \pm 1.24$ $75.32 \pm 0.83$ $73.08 \pm 4.35$ $20.16 \pm 13.44$	$74.92 \pm 0.80 74.91 \pm 0.97 72.35 \pm 8.86 20.06 \pm 15.76$
CIFAR-10	Vit-Small-192	384 256 128 64	$83.30 \pm 0.29$	$82.99 \pm 0.11$ $83.07 \pm 0.42$ $69.09 \pm 25.20$ $43.13 \pm 30.28$	$82.58 \pm 0.07  82.51 \pm 0.55  41.15 \pm 25.26  12.67 \pm 7.11$
CIFAR-10	ResNet-18 <sup>†</sup>	1024 512 128 64	$94.45 \pm 0.14$	$93.55 \pm 0.13$ $93.49 \pm 0.19$ $57.98 \pm 34.02$ $29.92 \pm 19.79$	$93.47 \pm 0.20$ $93.43 \pm 0.64$ $47.55 \pm 37.46$ $21.93 \pm 19.86$

<sup>†</sup>ResNet-18 results use PCA-reduced weights.

Impact of Canonicalization: We observe a capacity-dependent behavior of DeepWeightFlow models with and without canonicalization: i) at lower capacity of the FM models, models trained on canonicalized neural network weights generate higher performing ensembles than the FM models trained on non-canonicalized data. However, as the capacity of the FM model increases, the performances of the ensembles of generated neural networks become similar. In general, FM models trained on canonicalized neural network weights approach the performance of the training set ("original" neural networks) with lower capacity. Moreover, when flow model parameters are limited, models trained on canonicalized data generate neural networks with observably lower variance in accuracy compared to non-canonicalized counterparts. In Table 4 we show the performance of DeepWeightFlow with and without canonicalization.

**Robustness Across Initialization Schemes:** To evaluate generalization capability, we conducted extensive robustness testing using MLP models trained on the Iris dataset with diverse initialization strategies (Kaiming (He et al., 2015), Xavier (Glorot & Bengio, 2010), Kaiming weights and zero for biases, normal, and uniform distributions). *Training a single flow model on this heterogeneous collection (100 models total: 20 seeds* × 5 *initialization types) successfully generated novel weights achieving high test accuracy, demonstrating the framework's ability to learn from and generate weights across different initialization regimes.* All other experiments maintained consistency by using Kaiming initialization with varied random seeds.

#### 5.2 Transfer Learning on Unseen Datasets

Our generated models can be effectively used for transfer learning (Nava et al., 2023; Zhang et al., 2024b) across unseen datasets. In our experiments, we trained DeepWeightFlow on ResNet-18 models for the CIFAR-10 dataset using PCA, generated 5 models, and recalibrated their batch nor-

 $<sup>*</sup>d_h$ : flow hidden dimension

malization running mean and variance on a small subset of CIFAR-10. These models were then evaluated under zero-shot and finetuning settings on STL-10 and SVHN datasets. The results are presented in Table 5. DeepWeightFlow-generated models consistently outperformed state-of-the-art FM models such as FloWN (Saragih et al., 2025b) in both zero-shot and finetuning evaluations. Furthermore, they significantly outperformed randomly initialized models, proving the effectiveness of the method.

Table 5: Zero-shot (Epoch 0) and fine-tuning performance comparison of complete ResNet-18 CIFAR-10 classifiers generated by DeepWeightFlow, FLoWN, and RandomInit. RandomInit refers to fresh Kaiming-He initialization. Best results for each fine-tuning checkpoint are in bold. Saragih et al. (2025b)

Epoch	Model	Method	STL-10	SVHN
0	FLoWN	RandomInit Generated	$\begin{array}{c} 10.00 \pm 0.00 \\ 35.16 \pm 1.24 \end{array}$	$10.00 \pm 0.00$ <b>17.99</b> $\pm$ <b>0.82</b>
	DeepWeightFlow	RandomInit Generated	$11.18 \pm 1.48$ $48.32 \pm 0.34$	$8.01 \pm 1.41$ $11.57 \pm 0.49$
1	FLoWN	RandomInit Generated	$18.94 \pm 0.09$ $36.15 \pm 1.14$	$19.50 \pm 0.03$ $68.64 \pm 7.07$
	DeepWeightFlow	RandomInit Generated	$38.28 \pm 1.07$ <b>79.69</b> $\pm$ <b>1.08</b>	$84.07 \pm 1.76$ <b>91.66</b> $\pm$ <b>0.79</b>
5	FLoWN	RandomInit Generated	$28.24 \pm 0.01$ $37.43 \pm 1.19$	$39.59 \pm 10.0$ $77.36 \pm 1.07$
	DeepWeightFlow	RandomInit Generated	$51.35 \pm 0.51$ <b>84.63</b> $\pm$ <b>0.17</b>	$93.82 \pm 0.16$ $95.85 \pm 0.09$

#### 5.3 DIVERSITY OF GENERATED MODELS

To evaluate the DeepWeightFlow models' generative capabilities, we compute the maximum IoU (mIoU) between the generated neural networks and the neural networks in the training set (referred to as the "original" neural networks). The mIoU is constructed from the intersection over union of the wrong predictions made by the neural networks (Wang et al., 2024). It is defined as  $IoU = |P_1^{wrong} \cap P_2^{wrong}|/|P_1^{wrong} \cup P_2^{wrong}|$ , where  $P_1$  comes from the set being compared (such as from the generated set) and  $P_2$  comes from a reference set (such as the set of original neural networks). We disregard the IoU of a neural network with itself as it is trivially 1. The mIoU measure scales from complete dissimilarity at 0 to complete similarity at 1.

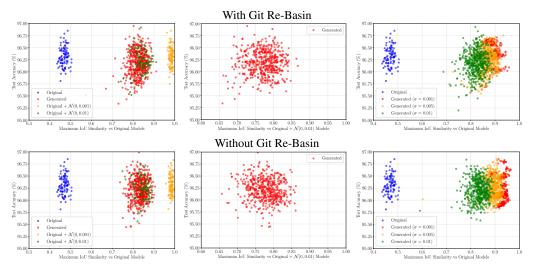


Figure 2: Maximum IoU vs test set accuracy for MNIST classifying MLPs. Lower maximum IoU implies greater diversity between the neural network weights. The left-most panels are generated and original neural networks (from the DeepWeightFlow training set) with different scales of Gaussian noise added with respect to the original neural networks. The middle panels show that the generated neural networks and the original neural networks with noise added, which overlap in the left-most panels, are concretely different. The right-most panels contain the original and generated neural networks with different source distributions. All panels include 500 models generated by DeepWeightFlow.

In Figure 2, we compare the original neural networks with the generated ones, with noise added to the weights of the original neural networks, and with neural networks generated with different FM source distributions. The upper row compares the cases for the FM models trained with re-basin, and the lower panels, without. In the left-most panels, we see that i) the original networks are quite diverse from each other, as evident from the blue cloud. This is the case as, unlike several previous works, we do not use checkpoints from the training of a single neural network as the training set of the DeepWeightFlow model. Instead, we generate independent sets of neural networks from random initializations. ii) The yellow and green clouds show that adding progressively increasing Gaussian noise to the original networks makes them progressively diverse from the original networks as expected (< 1). iii) The red cloud representing the generated networks shows diversity from the original set but seems to overlap with the green set, which represents the set created by adding noise sampled from  $\mathcal{N}(0,0.01)$  to the original neural network weights. From the middle panels in Figure 2, we see that the red cloud representing the generated neural networks is sufficiently diverse from the original ones with added noise sampled from  $\mathcal{N}(0, 0.01)$ . This gives us confidence that the generated neural networks are, indeed, not the same as the original networks with noise added to the weights. Lastly, the right-most panels show how diverse the generated neural networks are when generated with different source distributions. Hence, DeepWeightFlow is capable of generating a diverse set of neural networks while maintaining the accuracy of the task. In Appendix J, we provide the numerical estimates of mIoU, the Jensen-Shannon, Wasserstein, and Nearest Neighbors (NN) distances between generated and original neural networks.

#### 5.4 Training and Sampling Efficiency

DeepWeightFlow is significantly faster to train and generate neural network weights when compared to diffusion models in complete neural network weights generation. DeepWeightFlow takes up to  $\mathcal{O}(10)$  minutes to train for most neural network architectures with up to  $\mathcal{O}(10M)$  parameters as compared to the several hours that it takes to train RPG (Wang et al., 2025). For generating neural network weights, DeepWeightFlow takes a few minutes to generate hundreds of neural networks compared to the hours it takes to generate a comparable ensemble using RPG. Yet, DeepWeightFlow generates ensembles of neural networks that have comparable outcomes for ResNet-18s and ViTs. This is primarily because RPG is a diffusion model, whereas DeepWeightFlow is based on FM using a simple MLP implementation. A detailed comparison of training and generation efficiency can be found in Table 12 in Appendix I.

#### 6 Conclusion

In this work, we introduce DeepWeightFlow, a generative model for neural network weights that performs FM *directly* in weight space, unconditioned by dataset characteristics, task descriptions, or architectural specifications, and avoiding nonlinear dimensionality reduction. Through extensive experiments, we show that DeepWeightFlow generates diverse neural network weights for a variety of architectures (MLP, ResNet, ViT) that show excellent performance on vision and tabular classification tasks. We provide empirical evidence that canonicalizing the training data facilitates the generation of larger networks but is of limited usefulness for moderate-dimensional weights.

Scalability can be a concern when working directly in weight space. Even though our experiments show that DeepWeightFlow can be combined with simple linear dimensionality reduction techniques, such as PCA, to mitigate restrictions on model size, exploring scalability to very large models as well as the compatibility of DeepWeightFlow with model distillation, low-rank approximations, or sparsity remains future work. While we empirically investigate the role of canonicalization for weight generation, some open questions about the relative merits of canonicalization, equivariant architecture design, and data augmentation for learning in deep weight spaces remain.

#### REPRODUCIBILITY STATEMENT

The architectural details along with the hyperparameters used to generate the data have been provided in the main text and the appendix (Table 11, Table 10, and Table 9). The dataset will be made available on request and/or uploaded to a data repository. The code necessary to reproduce the results is in https://github.com/anonymousacademicc/DeepWeightFlow-ICLR.

#### REFERENCES

- Samuel Ainsworth, Jonathan Hayase, and Siddhartha Srinivasa. Git re-basin: Merging models modulo permutation symmetries. In *The Eleventh International Conference on Learning Representations*, 2023. URL https://openreview.net/forum?id=CQsmMYmlP5T.
- Jimmy Lei Ba, Jamie Ryan Kiros, and Geoffrey E. Hinton. Layer normalization, 2016. URL https://arxiv.org/abs/1607.06450.
- Chaoqi Chen, Weiping Xie, Wenbing Huang, Yu Rong, Xinghao Ding, Yue Huang, Tingyang Xu, and Junzhou Huang. Progressive feature alignment for unsupervised domain adaptation. In 2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), pp. 627–636, 2019. URL https://ieeexplore.ieee.org/document/8953748.
- Xiangning Chen. Advancing Automated Machine Learning: Neural Architectures and Optimization Algorithms. PhD thesis, University of California, Los Angeles, United States California, 2023. URL https://www.proquest.com/docview/2899619104/abstract/8CAD6EC2664A464CPQ/1.
- Adam Coates, Andrew Ng, and Honglak Lee. An Analysis of Single-Layer Networks in Unsupervised Feature Learning. In *Proceedings of the Fourteenth International Conference on Artificial Intelligence and Statistics*, pp. 215–223. JMLR Workshop and Conference Proceedings, June 2011. URL https://proceedings.mlr.press/v15/coates11a.html.
- Laurent Dinh, Razvan Pascanu, Samy Bengio, and Yoshua Bengio. Sharp minima can generalize for deep nets. In *Proceedings of the 34th International Conference on Machine Learning*, volume 70 of *Proceedings of Machine Learning Research*, pp. 1019–1028. PMLR, 2017. URL https://proceedings.mlr.press/v70/dinh17b.html.
- Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale. In *International Conference on Learning Representations*, 2021. URL https://openreview.net/forum?id=YicbFdNTTy.
- Rahim Entezari, Hanie Sedghi, Olga Saukh, and Behnam Neyshabur. The role of permutation invariance in linear mode connectivity of neural networks. In *International Conference on Learning Representations*, 2022. URL https://openreview.net/forum?id=dNigytemkL.
- Ege Erdogan. Geometric flow models over neural network weights, 2025. URL https://arxiv.org/abs/2504.03710.
- R. A. Fisher. The used of multiple measurements in taxonomic problems. *Annals of Eugenics*, 7(2): 179–188, 1936. URL https://onlinelibrary.wiley.com/doi/abs/10.1111/j. 1469–1809.1936.tb02137.x.
- Jonathan Frankle and Michael Carbin. The lottery ticket hypothesis: Finding sparse, trainable neural networks. In *International Conference on Learning Representations*, 2019. URL https://openreview.net/forum?id=rJl-b3RcF7.
- Xavier Glorot and Yoshua Bengio. Understanding the difficulty of training deep feedforward neural networks. In *Proceedings of the 13th International Conference on Artificial Intelligence and Statistics (AISTATS)*, pp. 249–256, 2010. URL http://proceedings.mlr.press/v9/glorot10a.html.
- Ian J. Goodfellow, Yaroslav Bulatov, Julian Ibarz, Sacha Arnoud, and Vinay Shet. Multi-digit number recognition from street view imagery using deep convolutional neural networks. In *Proceedings of the 2013 International Conference on Machine Learning (ICML)*, 2013. URL https://arxiv.org/abs/1312.6082.
- Ian J. Goodfellow, Jonathon Shlens, and Christian Szegedy. Explaining and harnessing adversarial examples. *arXiv*, 2015. URL https://arxiv.org/abs/1412.6572.

541

543

544

546

547

548

549

550

551

552

553

554

556

557

558

559

561

563

565

566

567 568

569

570 571

572

573

574

575 576

577

578

579

580 581

582

583

584 585

586

588

589

590

592

Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Peiyi Wang, Qihao Zhu, Runxin Xu, Ruoyu Zhang, Shirong Ma, Xiao Bi, Xiaokang Zhang, Xingkai Yu, Yu Wu, Z. F. Wu, Zhibin Gou, Zhihong Shao, Zhuoshu Li, Ziyi Gao, Aixin Liu, Bing Xue, Bingxuan Wang, Bochao Wu, Bei Feng, Chengda Lu, Chenggang Zhao, Chengqi Deng, Chong Ruan, Damai Dai, Deli Chen, Dongjie Ji, Erhang Li, Fangyun Lin, Fucong Dai, Fuli Luo, Guangbo Hao, Guanting Chen, Guowei Li, H. Zhang, Hanwei Xu, Honghui Ding, Huazuo Gao, Hui Qu, Hui Li, Jianzhong Guo, Jiashi Li, Jingchang Chen, Jingyang Yuan, Jinhao Tu, Junjie Qiu, Junlong Li, J. L. Cai, Jiaqi Ni, Jian Liang, Jin Chen, Kai Dong, Kai Hu, Kaichao You, Kaige Gao, Kang Guan, Kexin Huang, Kuai Yu, Lean Wang, Lecong Zhang, Liang Zhao, Litong Wang, Liyue Zhang, Lei Xu, Leyi Xia, Mingchuan Zhang, Minghua Zhang, Minghui Tang, Mingxu Zhou, Meng Li, Miaojun Wang, Mingming Li, Ning Tian, Panpan Huang, Peng Zhang, Qiancheng Wang, Qinyu Chen, Qiushi Du, Ruiqi Ge, Ruisong Zhang, Ruizhe Pan, Runji Wang, R. J. Chen, R. L. Jin, Ruyi Chen, Shanghao Lu, Shangyan Zhou, Shanhuang Chen, Shengfeng Ye, Shiyu Wang, Shuiping Yu, Shunfeng Zhou, Shuting Pan, S. S. Li, Shuang Zhou, Shaoqing Wu, Tao Yun, Tian Pei, Tianyu Sun, T. Wang, Wangding Zeng, Wen Liu, Wenfeng Liang, Wenjun Gao, Wenqin Yu, Wentao Zhang, W. L. Xiao, Wei An, Xiaodong Liu, Xiaohan Wang, Xiaokang Chen, Xiaotao Nie, Xin Cheng, Xin Liu, Xin Xie, Xingchao Liu, Xinyu Yang, Xinyuan Li, Xuecheng Su, Xuheng Lin, X. Q. Li, Xiangyue Jin, Xiaojin Shen, Xiaosha Chen, Xiaowen Sun, Xiaoxiang Wang, Xinnan Song, Xinyi Zhou, Xianzu Wang, Xinxia Shan, Y. K. Li, Y. Q. Wang, Y. X. Wei, Yang Zhang, Yanhong Xu, Yao Li, Yao Zhao, Yaofeng Sun, Yaohui Wang, Yi Yu, Yichao Zhang, Yifan Shi, Yiliang Xiong, Ying He, Yishi Piao, Yisong Wang, Yixuan Tan, Yiyang Ma, Yiyuan Liu, Yongqiang Guo, Yuan Ou, Yuduan Wang, Yue Gong, Yuheng Zou, Yujia He, Yunfan Xiong, Yuxiang Luo, Yuxiang You, Yuxuan Liu, Yuyang Zhou, Y. X. Zhu, Yanping Huang, Yaohui Li, Yi Zheng, Yuchen Zhu, Yunxian Ma, Ying Tang, Yukun Zha, Yuting Yan, Z. Z. Ren, Zehui Ren, Zhangli Sha, Zhe Fu, Zhean Xu, Zhenda Xie, Zhengyan Zhang, Zhewen Hao, Zhicheng Ma, Zhigang Yan, Zhiyu Wu, Zihui Gu, Zijia Zhu, Zijun Liu, Zilin Li, Ziwei Xie, Ziyang Song, Zizheng Pan, Zhen Huang, Zhipeng Xu, Zhongyu Zhang, and Zhen Zhang. DeepSeek-R1 incentivizes reasoning in LLMs through reinforcement learning. Nature, 645 (8081):633–638, September 2025. ISSN 1476-4687. doi: 10.1038/s41586-025-09422-z. URL https://www.nature.com/articles/s41586-025-09422-z.

David Ha, Andrew M. Dai, and Quoc V. Le. Hypernetworks. In *International Conference on Learning Representations*, 2017. URL https://openreview.net/forum?id=rkpACellx.

Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Delving deep into rectifiers: Surpassing human-level performance on imagenet classification. In *Proceedings of the 2015 IEEE International Conference on Computer Vision (ICCV)*, pp. 1026-1034, 2015. URL https://www.cv-foundation.org/openaccess/content\_iccv\_2015/html/He\_Delving\_Deep\_into\_ICCV\_2015\_paper.html.

Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep Residual Learning for Image Recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 770–778, 2016. URL https://openaccess.thecvf.com/content\_cvpr\_2016/html/He\_Deep\_Residual\_Learning\_CVPR\_2016\_paper.html.

Robert Hecht-Nielsen. On the algebraic structure of feedforward network weight spaces. In Rolf Eckmiller (ed.), *Advanced Neural Computers*, pp. 129–135. North-Holland, Amsterdam, 1990. ISBN 978-0-444-88400-8. URL https://www.sciencedirect.com/science/article/pii/B9780444884008500194.

Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In Francis Bach and David Blei (eds.), *Proceedings of the 32nd International Conference on Machine Learning*, volume 37 of *Proceedings of Machine Learning Research*, pp. 448–456, Lille, France, 07–09 Jul 2015. PMLR. URL https://proceedings.mlr.press/v37/ioffe15.html.

Pavel Izmailov, Dmitrii Podoprikhin, Timur Garipov, Dmitry Vetrov, and Andrew Gordon Wilson. Averaging weights leads to wider optima and better generalization. In *Proceedings of the 34th Conference on Uncertainty in Artificial Intelligence (UAI)*, pp. 876–885. AUAI Press, 2018. URL https://arxiv.org/abs/1803.05407.

- R. Jonker and A. Volgenant. A shortest augmenting path algorithm for dense and sparse linear assignment problems. *Computing*, 38(4):325–340, November 1987. ISSN 0010-485X. URL https://doi.org/10.1007/BF02278710.
- Keller Jordan, Hanie Sedghi, Olga Saukh, Rahim Entezari, and Behnam Neyshabur. Repair: Renormalizing permuted activations for interpolation repair. *arXiv* preprint arXiv:2211.08403, 2022. doi: 10.48550/arXiv.2211.08403. URL https://arxiv.org/abs/2211.08403.
- Bobak Kiani, Jason Wang, and Melanie Weber. Hardness of learning neural networks under the manifold hypothesis. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024. URL https://openreview.net/forum?id=dkkgKzMni7.
- Boris Knyazev, Michal Drozdzal, Graham W. Taylor, and Adriana Romero. Parameter prediction for unseen deep architectures. In A. Beygelzimer, Y. Dauphin, P. Liang, and J. Wortman Vaughan (eds.), *Advances in Neural Information Processing Systems*, 2021. URL https://openreview.net/forum?id=vqHak8NLk25.
- Alex Krizhevsky, Vinod Nair, and Geoffrey Hinton. Learning multiple layers of features from tiny images. Technical report, University of Toronto, 2009. URL http://www.cs.toronto.edu/kriz/cifar.html.
- Balaji Lakshminarayanan, Alexander Pritzel, and Charles Blundell. Simple and scalable predictive uncertainty estimation using deep ensembles. In *Proceedings of the 31st International Conference on Neural Information Processing Systems*, NIPS'17, pp. 6405–6416, Red Hook, NY, USA, 2017. Curran Associates Inc. ISBN 9781510860964. URL https://dl.acm.org/doi/10.5555/3295222.3295387.
- Y. Lecun, L. Bottou, Y. Bengio, and P. Haffner. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998. URL https://ieeexplore.ieee.org/document/726791.
- Derek Lim, Theo Putterman, Robin Walters, Haggai Maron, and Stefanie Jegelka. The empirical impact of neural parameter symmetries, or lack thereof. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024. URL https://openreview.net/forum?id=pCVxYw6FKg.
- Yaron Lipman, Ricky T. Q. Chen, Heli Ben-Hamu, Maximilian Nickel, and Matthew Le. Flow matching for generative modeling. In *The Eleventh International Conference on Learning Representations*, 2023. URL https://openreview.net/forum?id=PqvMRDCJT9t.
- Ziyin Liu. Symmetry leads to structured constraint of learning, 2023. URL https://arxiv.org/abs/2309.16932.
- Wesley J Maddox, Timur Garipov, Pavel Izmailov, Dmitry Vetrov, and Andrew Gordon Wilson. A simple baseline for bayesian uncertainty in deep learning. *arXiv preprint arXiv:1902.02476*, 2019. URL https://arxiv.org/abs/1902.02476.
- Aleksander Madry, Aleksandar Makelov, Ludwig Schmidt, Dimitris Tsipras, and Adrian Vladu. Towards deep learning models resistant to adversarial attacks. In *International Conference on Learning Representations*, 2018. URL https://openreview.net/forum?id=rJzIBfZAb.
- Kevin Meng, David Bau, Alex J Andonian, and Yonatan Belinkov. Locating and editing factual associations in GPT. In Alice H. Oh, Alekh Agarwal, Danielle Belgrave, and Kyunghyun Cho (eds.), *Advances in Neural Information Processing Systems*, 2022. URL https://openreview.net/forum?id=-h6WAS6eE4.
- Eric Mitchell, Charles Lin, Antoine Bosselut, Chelsea Finn, and Christopher D Manning. Fast model editing at scale. In *International Conference on Learning Representations*, 2022. URL https://openreview.net/forum?id=0DcZxeWfOPt.
- Milad Nasr, Reza Shokri, and Amir Houmansadr. Comprehensive privacy analysis of deep learning: Passive and active white-box inference attacks against centralized and federated learning. In 2019 IEEE symposium on security and privacy (SP), pp. 739–753. IEEE, 2019. URL https://www.computer.org/csdl/proceedings-article/sp/2019/666000a739/1dlwhtj4r70.

- Elvis Nava, Seijin Kobayashi, Yifei Yin, Robert K. Katzschmann, and Benjamin F. Grewe. Meta-learning via classifier(-free) diffusion guidance. *Transactions on Machine Learning Research*, 4: 1–20, 2023. URL https://openreview.net/forum?id=lirVjE7A3w.
- Aviv Navon, Aviv Shamsian, Idan Achituve, Ethan Fetaya, Gal Chechik, and Haggai Maron. Equivariant architectures for learning in deep weight spaces. In *Proceedings of the 40th International Conference on Machine Learning*, ICML'23. JMLR.org, 2023. URL https://dl.acm.org/doi/10.5555/3618408.3619481.
- Behnam Neyshabur, Ruslan Salakhutdinov, and Nati Srebro. Path-sgd: Path-normalized optimization in deep neural networks. In Advances in Neural Information Processing Systems, volume 28, 2015a. URL https://papers.nips.cc/paper/5797-path-sgd-path-normalized-optimization-in-deep-neural-networks.pdf.
- Behnam Neyshabur, Ryota Tomioka, and Nathan Srebro. Norm-based capacity control in neural networks. In *Proceedings of the 28th Conference on Learning Theory*, volume 40 of *Proceedings of Machine Learning Research*, pp. 1376–1401. PMLR, 2015b. URL https://proceedings.mlr.press/v40/Neyshabur15.html.
- William Peebles, Ilija Radosavovic, Tim Brooks, Alexei A. Efros, and Jitendra Malik. Learning to learn with generative models of neural network checkpoints, 2022. URL https://arxiv.org/abs/2209.12892.
- Fabrizio Pittorino, Antonio Ferraro, Gabriele Perugini, Christoph Feinauer, Carlo Baldassi, and Riccardo Zecchina. Deep networks on toroids: Removing symmetries reveals the structure of flat regions in the landscape geometry. In *Proceedings of the 39th International Conference on Machine Learning*, pp. 17759–17781. PMLR, 2022. URL https://proceedings.mlr.press/v162/pittorino22a.html.
- Neale Ratzlaff and Li Fuxin. HyperGAN: A generative model for diverse, performant neural networks. In Kamalika Chaudhuri and Ruslan Salakhutdinov (eds.), *Proceedings of the 36th International Conference on Machine Learning*, volume 97 of *Proceedings of Machine Learning Research*, pp. 5361–5369. PMLR, 09–15 Jun 2019. URL https://proceedings.mlr.press/v97/ratzlaff19a.html.
- Filippo Rinaldi, Giacomo Capitani, Lorenzo Bonicelli, Donato Crisostomi, Federico Bolelli, ELISA FICARRA, Emanuele Rodolà, Simone Calderara, and Angelo Porrello. Update your transformer to the latest release: Re-basin of task vectors. In *Forty-second International Conference on Machine Learning*, 2025. URL https://openreview.net/forum?id=sHvImzN9pL.
- Daniel Saragih, Deyu Cao, and Tejas Balaji. Flows and diffusions on the neural manifold, 2025a. URL https://arxiv.org/abs/2507.10623.
- Daniel Saragih, Deyu Cao, Tejas Balaji, and Ashwin Santhosh. Flow to learn: Flow matching on neural network parameters. In *Workshop on Neural Network Weights as a New Data Modality*, 2025b. URL https://openreview.net/forum?id=r0ynTstq3c.
- Konstantin Schürholt, Boris Knyazev, Xavier Giró i Nieto, and Damian Borth. Hyper-representations as generative models: Sampling unseen neural network weights. In Alice H. Oh, Alekh Agarwal, Danielle Belgrave, and Kyunghyun Cho (eds.), Advances in Neural Information Processing Systems, 2022. URL https://openreview.net/forum?id=uyEYNg2HHFQ.
- Konstantin Schürholt, Michael W. Mahoney, and Damian Borth. Towards scalable and versatile weight space learning. In *Forty-first International Conference on Machine Learning*, 2024. URL https://openreview.net/forum?id=ug2uoAZ9c2.
- Aviv Shamsian, David Zhang, Aviv Navon, Yan Zhang, Miltiadis Kofinas, Idan Achituve, Riccardo Valperga, Gertjan Burghouts, Efstratios Gavves, Cees Snoek, Ethan Fetaya, Gal Chechik, and Haggai Maron. Data Augmentations in Deep Weight Spaces. In *NeurIPS 2023 Workshop on Symmetry and Geometry in Neural Representations*, November 2023. URL https://openreview.net/forum?id=jdT7PuqdSt.

```
Aviv Shamsian, Aviv Navon, David W. Zhang, Yan Zhang, Ethan Fetaya, Gal Chechik, and Haggai Maron. Improved generalization of weight space networks via augmentations. In Proceedings of the 41st International Conference on Machine Learning, volume 235 of ICML'24, pp. 44378–44393, Vienna, Austria, July 2024. JMLR.org. URL https://dl.acm.org/doi/abs/10.5555/3692070.3693876.
```

- Gil Shomron and Uri Weiser. Post-training batchnorm recalibration, 2020. URL https://arxiv.org/abs/2010.05625.
- Bedionita Soro, Bruno Andreis, Hayeon Lee, Wonyong Jeong, Song Chong, Frank Hutter, and Sung Ju Hwang. Diffusion-based neural network weights generation. In *The Thirteenth International Conference on Learning Representations*, 2025. URL https://openreview.net/forum?id=j8WHjM9aMm.
- Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurelien Rodriguez, Armand Joulin, Edouard Grave, and Guillaume Lample. LLaMA: Open and Efficient Foundation Language Models, February 2023. URL http://arxiv.org/abs/2302.13971.
- Florian Tramer, Andreas Terzis, Thomas Steinke, Shuang Song, Matthew Jagielski, and Nicholas Carlini. Debugging differential privacy: A case study for privacy auditing. *arXiv*, 2022. URL https://arxiv.org/abs/2202.12219.
- Dequan Wang, Evan Shelhamer, Shaoteng Liu, Bruno Olshausen, and Trevor Darrell. Tent: Fully test-time adaptation by entropy minimization. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pp. 4068–4078, 2021. URL https://arxiv.org/abs/2103.06905.
- Kai Wang, Zhaopan Xu, Yukun Zhou, Zelin Zang, Trevor Darrell, Zhuang Liu, and Yang You. Neural Network Diffusion, February 2024. URL http://arxiv.org/abs/2402.13144.
- Kai Wang, Dongwen Tang, Wangbo Zhao, Konstantin Schürholt, Zhangyang Wang, and Yang You. Recurrent diffusion for large-scale parameter generation, 2025. URL https://arxiv.org/abs/2501.11587.
- Mitchell Wortsman, Maxwell C Horton, Carlos Guestrin, Ali Farhadi, and Mohammad Rastegari. Learning neural network subspaces. In Marina Meila and Tong Zhang (eds.), *Proceedings of the 38th International Conference on Machine Learning*, volume 139 of *Proceedings of Machine Learning Research*, pp. 11217–11227. PMLR, 18–24 Jul 2021. URL https://proceedings.mlr.press/v139/wortsman21a.html.
- Mitchell Wortsman, Gabriel Ilharco, Samir Ya Gadre, Rebecca Roelofs, Raphael Gontijo-Lopes, Ari S Morcos, Hongseok Namkoong, Ali Farhadi, Yair Carmon, Simon Kornblith, and Ludwig Schmidt. Model soups: averaging weights of multiple fine-tuned models improves accuracy without increasing inference time. In Kamalika Chaudhuri, Stefanie Jegelka, Le Song, Csaba Szepesvari, Gang Niu, and Sivan Sabato (eds.), *Proceedings of the 39th International Conference on Machine Learning*, volume 162 of *Proceedings of Machine Learning Research*, pp. 23965–23998. PMLR, 17–23 Jul 2022. URL https://proceedings.mlr.press/v162/wortsman22a.html.
- Han Xiao, Kashif Rasul, and Roland Vollgraf. Fashion-mnist: a novel image dataset for benchmarking machine learning algorithms, 2017. URL https://arxiv.org/abs/1708.07747.
- Boya Zeng, Yida Yin, Zhiqiu Xu, and Zhuang Liu. Generative modeling of weights: Generalization or memorization?, 2025. URL https://arxiv.org/abs/2506.07998.
- Baoquan Zhang, Chuyao Luo, Demin Yu, Xutao Li, Huiwei Lin, Yunming Ye, and Bowen Zhang. Metadiff: Meta-learning with conditional diffusion for few-shot learning. *Proceedings of the AAAI Conference on Artificial Intelligence*, 38(15):16687–16695, Mar. 2024a. URL https://ojs.aaai.org/index.php/AAAI/article/view/29608.
- Baoquan Zhang, Chuyao Luo, Demin Yu, Xutao Li, Huiwei Lin, Yunming Ye, and Bowen Zhang. Metadiff: Meta-learning with conditional diffusion for few-shot learning. In *Proceedings of the AAAI conference on artificial intelligence*, volume 38, pp. 16687–16695, 2024b.

Bo Zhao, Robert M. Gower, Robin Walters, and Rose Yu. Improving convergence and generalization using parameter symmetries. In *International Conference on Learning Representations*, 2024. URL https://openreview.net/forum?id=L0r0GphlIL.

Bo Zhao, Robin Walters, and Rose Yu. Symmetry in Neural Network Parameter Spaces, June 2025. URL http://arxiv.org/abs/2506.13018.

## A GIT RE-BASIN

Git Re-Basin weight matching, formulated by Ainsworth et al. (2023), is a greedy permutation coordinate descent algorithm for moving a model's weights  $\theta_A$  into the same 'basin' in the loss landscape of the model class  $f_{\hat{\theta}}$  as a reference model's weights  $\theta_B$ .

This operation is applied here as a canonicalization step before weight flattening and the subsequent training of the DeepWeightFlow models. The procedure reduces the space of the task from  $\mathbb{R}^{\theta}$  to a quotient space of  $\mathbb{R}^{\theta}$  modulo permutation symmetry.

Applying this across the model layers constructs a transformed model  $\theta'$  by

$$W'_{\ell} = PW_{\ell}, \ b'_{\ell} = Pb_{\ell}, \ W'_{\ell+1} = W_{\ell+1}P^{T}$$
(3)

The 'distance' between two permutations is therefore a Frobenius inner product of  $P_\ell W_\ell^A$  and  $W_\ell^B$ , written as  $\langle A,B\rangle=\sum_{i,j}A_{i,j}B_{i,j}$  for real-valued matrices A and B. Accounting for the transforms outlined above, the process of matching the permutations across the stack of layers becomes,

$$\underset{\pi = \{P_{\ell}\}_{L}^{L}}{\arg \max} \sum_{n=1}^{L} \left\langle W_{i}^{B}, P_{i} W_{i}^{A} P_{i-1}^{T} \right\rangle \text{ with } P_{0}^{T} = I \tag{4}$$

This formulation presents a Symmetric Orthogonal Bilinear Assignment Problem (SOBLAP), which is NP-hard. However, when relaxed to focus on a single permutation  $P_{\ell}$  at a time - *ceteris paribus*, the problem simplifies to a series of Linear Assignment Problems (LAPs) of the form below (Ainsworth et al., 2023; Zhao et al., 2025; Rinaldi et al., 2025). These LAPs can be solved in polynomial time by methods like the Hungarian Algorithm (Jonker & Volgenant, 1987).

$$\arg\max_{P_{\ell}} \left\langle W_{\ell}^{B}, P_{\ell} W_{\ell}^{A} P_{\ell-1}^{T} \right\rangle + \left\langle W_{\ell+1}^{B}, P_{\ell+1} W_{\ell+1}^{A} P_{\ell}^{T} \right\rangle \tag{5}$$

The product of this process is a permutation  $\pi'$  of model A's weights into the same basin in  $f_{\theta}$ 's loss landscape as model B with exact functional equivalence ( $f_{\theta_A} = f_{\pi'(\theta_A)}$ ). However, sequences of LAPs are understood to be coarse approximations of SOBLAPs and, as such, strong conclusions cannot be drawn about the optimality of  $\pi'$  (Rinaldi et al., 2025; Ainsworth et al., 2023).

#### **B** Transfusion

We canonicalize a collection of Vision Transformers (ViTs) using the method of Rinaldi et al. (2025), which introduces a structured alignment procedure for multi-head attention transformer weights (Rinaldi et al., 2025).

The core difficulty in transformers arises from multi-head attention and residual connections: Naive global permutations either mix information across heads or break functional equivalence in residual branches (Zhao et al., 2025). To address this, the method applies a *two-level permutation scheme*:

 Inter-Head Alignment: For each multi-head attention layer, attention heads from different checkpoints are first matched. This is done by comparing the singular value spectra of their projection matrices, which are invariant under row and column permutations, and then solving the resulting assignment problem with the Hungarian algorithm. This step ensures that corresponding heads are correctly paired across models.

For a sub matrix representing a single attention head in model A,  $h_i^A = [\tilde{W}]_i^A \in \mathbb{R}^{k \times m}$ , where k is the key value dimension and m is the attention embedding dimension, apply

singular value decomposition  $(A = U\Sigma V^T)$  to access the spectral projection matricies  $\Sigma$ , which are invariant to row and column permutations. For every head in a layer of model A, construct a distance,  $d_i, j = ||\Sigma_i - \Sigma_j||$ . These distances can be constructed for q, k, and v for each head and combined linearly  $D_{i,j} = d^q_{i,j} + d^k_{i,j} + d^v_{i,j}$  with  $D_{i,j} \in \mathbb{R}^{H \times H}$  (H is the number of heads). Therefore the optimal pairing of heads for model A and B is (Rinaldi et al., 2025),

$$P_{\text{inter head}} = \underset{P \in S_H}{\arg\min} \sum D_{i,P[i]}$$
 (6)

2. **Intra-Head Alignment:** Once heads are paired, the method refines the alignment by permuting rows and columns *within* each head independently, again solved via assignment on pairwise similarity scores. Restricting permutations within heads preserves head isolation and guarantees that residual connections remain valid after alignment.

After matching the heads of A to B the goal aligns closely with Git-ReBasin (Ainsworth et al., 2023) - to reorder  $h_{P[i]}^A$  such that the Frobenius inner product is maximized between H sub portions (Rinaldi et al., 2025),

$$P_{\text{intra head}}^{(i)} = \arg\max\langle h_i^B, Ph_{P[i]}^A \rangle \tag{7}$$

By iterating these two stages across all transformer layers, the procedure yields a canonicalized parameterization in which weights are aligned up to permutation symmetries. The goal is to permute units in such a way that two weight sets  $\theta_A$  and  $\theta_B$  become functionally comparable, reducing the effective size of the weight space that the FM encounters Rinaldi et al. (2025). This is similar to the case of Git Re-Basin (Ainsworth et al., 2023) for canonicalization.

## C RECALIBRATION OF BATCH NORMALIZATION WEIGHTS

Given a generated neural network with randomly initialized or flow-matched weights, the batch normalization layers contain statistics that may not match the actual data distribution. Naively interpolating weights of trained networks can lead to variance collapse (Jordan et al., 2022; Ainsworth et al., 2023), where the per-channel activation variances shrink drastically, breaking normalization and degrading performance. The recalibration process computes proper running statistics using the target dataset(Izmailov et al., 2018; Maddox et al., 2019; Shomron & Weiser, 2020; Wang et al., 2021).

We include these statistics parameters of batch normalization layers in the PermutationSpec of Git Re-Basin, a config that defines the permutation ordering across layers for weight matching, so that these statistics are also permuted and correctly maintained, ensuring that the permuted networks retain the same weights and accuracy as the original network.

#### C.1 STANDARD BATCH NORMALIZATION

For a feature map  $\mathbf{x} \in \mathbb{R}^{N \times C \times H \times W}$  where N is batch size, C is channels, and H, W are spatial dimensions:

$$\mu_c = \frac{1}{NHW} \sum_{n=1}^{N} \sum_{h=1}^{H} \sum_{w=1}^{W} x_{n,c,h,w}$$
(8)

$$\sigma_c^2 = \frac{1}{NHW} \sum_{n=1}^N \sum_{h=1}^H \sum_{w=1}^W (x_{n,c,h,w} - \mu_c)^2$$
 (9)

$$\hat{x}_{n,c,h,w} = \frac{x_{n,c,h,w} - \mu_c}{\sqrt{\sigma_c^2 + \epsilon}} \tag{10}$$

$$y_{n,c,h,w} = \gamma_c \hat{x}_{n,c,h,w} + \beta_c \tag{11}$$

where  $\gamma_c$  and  $\beta_c$  are learnable scale and shift parameters, and  $\epsilon$  is a small constant for numerical stability. During training, BatchNorm (Ioffe & Szegedy, 2015) maintains running statistics using an exponential moving average:

$$\bar{\mu}_c^{(t)} = (1 - \alpha)\bar{\mu}_c^{(t-1)} + \alpha\mu_c^{(t)}$$
 (12)

867 868  $\bar{\sigma}_{c}^{2(t)} = (1 - \alpha)\bar{\sigma}_{c}^{2(t-1)} + \alpha\sigma_{c}^{2(t)}$ (13)

where  $\alpha$  is the momentum parameter, typically 0.1, and t denotes the time step.

870

871

872

873

874

875

876

878

879

880

883

885

886 887

889 890

891 892 893

894

## **Algorithm 1** Batch Normalization Recalibration

- 1: **Input:** Calibration dataset  $\mathcal{D}$  (e.g., test dataset), batch size B
- 2: H and W denote the height and width of feature maps
- 3:  $x_{i,c,h,w}$  denotes the activation of sample i, channel c, at spatial position (h, w).
- 4: Initialize  $\bar{\mu}_c=0,\,\bar{\sigma}_c^2=1,\,n_c=0$  for all channels c
- 5: Disable exponential moving average (momentum) updates
- 6: Partition  $\mathcal{D}$  into mini-batch sequence  $\{\mathcal{B}_1, \mathcal{B}_2, \dots, \mathcal{B}_K\}$  where  $\bigcup_{k=1}^K \mathcal{B}_k = \mathcal{D}$
- 7: Define batch statistics for each  $\mathcal{B}_k$  and channel c:

$$\mu_c^{(k)} = \frac{1}{|\mathcal{B}_k|HW} \sum_{i \in \mathcal{B}_k} \sum_{h=1}^H \sum_{w=1}^W x_{i,c,h,w}$$
$$\sigma_c^{2(k)} = \frac{1}{|\mathcal{B}_k|HW} \sum_{i \in \mathcal{B}_k} \sum_{h=1}^H \sum_{w=1}^W (x_{i,c,h,w} - \mu_c^{(k)})^2$$

8: Compute running statistics where  $n_k = |\mathcal{B}_k|HW$  and  $n_c^{(k)} = n_c^{(k-1)} + n_k$ :

$$\begin{split} \bar{\mu}_c^{(k)} &= \frac{n_c^{(k-1)} \bar{\mu}_c^{(k-1)} + n_k \cdot \mu_c^{(k)}}{n_c^{(k)}} \\ \bar{\sigma}_c^{2(k)} &= \frac{n_c^{(k-1)} \bar{\sigma}_c^{2(k-1)} + n_k \cdot \sigma_c^{2(k)} + \frac{n_c^{(k-1)} n_k}{n_c^{(k)}} \left( \bar{\mu}_c^{(k-1)} - \mu_c^{(k)} \right)^2}{n_c^{(k)}} \end{split}$$

- 9: Final recalibrated statistics:  $\bar{\mu}_c = \bar{\mu}_c^{(K)}, \, \bar{\sigma}_c^2 = \bar{\sigma}_c^{2(K)}$  for all channels c
- 10: Restore exponential moving average updates (set momentum = 0.1)

895 897

## C.2 RECALIBRATION PROCESS

899 900 901

898

For generated networks, recompute running BatchNorm statistics:

902

- 1. **Reset**: Initialize running mean and variance for all channels, and set total sample count to
- 903
- 2. **Disable momentum**: Turn off exponential moving average updates.

904 905

3. Forward pass and incremental update: For each mini-batch in the calibration dataset:

906 907

• Compute the mean and variance of the batch for each channel.

908 909 • Update the running mean as a weighted average of the previous running mean and the batch mean.

910 911

 Update the running variance by combining the previous variance, the batch variance, and a correction for the shift in means. • Update the total sample count.

912

4. **Restore momentum**: Re-enable exponential moving average updates with the original momentum value.

914 915

916

917

The algorithm we use for recalibration of the batch normalization running statistics is provided in Algorithm 1. In Table 6 we show the results of recalibration on the generated neural networks. This clearly shows the importance of batch normalization, running statistics recalibration on the generation of neural networks that have batch normalization in their architecture.

Table 6: Comparing the impact of batch norm recalibration on complete ResNet-18 and 20s generated by DeepWeightFlow. Recalibrating batch normalization statistics on a small subset of target data significantly improves the accuracy of generated models.

Model	Git Re-Basin	Strategy	Mean $\pm$ Std (%)	Min (%)	Max (%)
ResNet-18	Yes	No Calibration Ref BN* Recalibrated	$10.00 \pm 0.00$ $19.06 \pm 9.68$ $93.05 \pm 4.42$	10.00 10.00 49.12	10.00 94.05 93.93
ResNet-18	No	No Calibration Ref BN Recalibrated	$10.00 \pm 0.00$ $10.28 \pm 1.24$ $93.49 \pm 0.21$	10.00 6.23 92.77	10.00 15.93 93.96
ResNet-20	Yes	No Calibration Ref BN Recalibrated	$14.36 \pm 3.10$ $17.88 \pm 4.66$ $74.57 \pm 0.84$	5.84 9.96 71.47	19.03 26.54 76.17
ResNet-20	No	No Calibration Ref BN Recalibrated	$12.64 \pm 2.22$ $10.23 \pm 0.79$ <b>75.21</b> $\pm$ <b>0.79</b>	8.12 8.04 72.06	18.19 14.92 76.52

<sup>\*</sup> Ref BN: Uses batch normalization statistics from reference model (seed 0)

## D FINETUNING MODELS FOR TRANSFER LEARNING ON UNSEEN DATASETS

We leverage ResNet-18 models trained and generated on the CIFAR-10 dataset to adapt to other unseen datasets, specifically STL-10 and SVHN (Table 5). We first evaluate the performance of the generated CIFAR-10 models on these datasets without any fine-tuning (Epoch 0). Subsequently, we fine-tune the models using the standard training set of the target dataset and evaluate them on the corresponding test set. Fine-tuning is performed for up to 5 epochs using the AdamW optimizer with a learning rate of  $1 \times 10^{-4}$ , weight decay of  $1 \times 10^{-4}$ , and a cosine learning rate scheduler with  $T_{\rm max} = epochs$  for smooth decay. We use a detach ratio of 0.4 (same as used by Saragih et al. (2025b)) and the cross-entropy loss is used as the objective function.

#### E PCA AS AN EFFECTIVE COMPRESSION STRATEGY

Table 7: Comparison of complete and PCA compressed weights generated by DeepWeightFlow. PCA enables stable generation at larger scales while improving accuracy for ResNet-20 and maintaining performance for Vit-Small-192. Best results are highlighted in bold.

		_	Original	Genera	ted Models
Model	Method $d_h$			With Re-basin	Without Re-basin
			Mean	Mean ± Std	Mean ± Std
ResNet-20 ResNet-20	Without PCA With PCA	512 512	$73.62 \pm 2.24$ $73.62 \pm 2.24$	75.07 ± 1.24 <b>75.96</b> ± <b>0.89</b>	74.92 ± 0.80 <b>75.97 ± 0.86</b>
Vit-Small-192 Vit-Small-192	Without PCA With PCA	384 1024	$83.30 \pm 0.29$ $83.30 \pm 0.29$	82.99 ± 0.11 83.08 ± 0.19	82. 58 ± 0.07 83.28 ± 0.01

In Table 7, we show the effects of using PCA to reduce the dimension of the neural network weight space. This is necessary as DeepWeightFlow cannot be trained on with the full rank of the larger neural networks, such as ResNet-18, due to memory constraints on a single GPU. Hence, we reduce dimensionality using PCA and decompress after generation. To test the validity of PCA, we trained the DeepWeightFlow models on ResNet-20 and ViT with and without using PCA as shown in Table 7. We observe that the accuracy and diversity of the neural networks (indicated by the standard deviation in the accuracy) are sufficiently representative of the original sample with or without PCA. This gives us confidence that much larger neural networks can be generated by DeepWeightFlow using PCA. We leave the complete implementation of this as future work.

Here we have performed incremental PCA that lets us perform PCA in chunks without loading all data into memory, but the math and essential foundation for it is exactly the same as standard PCA. Incremental PCA reduces the dimensionality of the generated weight matrices, we start with data of shape  $(n_{\text{samples}}, \text{flat\_dim})$ , incremental PCA projects it into a latent space of size  $(n_{\text{samples}}, \text{latent\_dim})$ , where we set latent\\_dim = 99. Since PCA orders components by explained variance and the rank of the data matrix is bounded by  $n_{\text{samples}} - 1$ , at most 99 meaningful directions

can exist for 100 samples we used. Therefore, using 99 principal components retains essentially all the variance of the dataset, while compressing the original high-dimensional representation into a very compact latent space.

## F NEED FOR CHOOSING THE RIGHT SOURCE DISTRIBUTION

The choice of source distribution for these generative models has a significant impact on the performance of the generated models. Table Table 8 highlights the importance of selecting a source distribution that aligns well with the target distributions to ensure reliable and high-quality weight generation.

Table 8: Evaluating the impact of various source distribution choices in FM mapping on the performance of complete weights generated by Deep-WeightFlow.

Model & Source Distribution	With Rebasin (%)	Without Rebasin (%)
Vit-Small-192 on CIFAR-10		
Original Accuracy	83.2	$29 \pm 0.29$
Gaussian(0, 0.01)	$78.31 \pm 10.99$	$76.69 \pm 14.37$
Gaussian(0, 0.001)	$\textbf{82.90} \pm \textbf{0.70}$	$82.40 \pm 5.29$
MLP on MNIST		
Original Accuracy	96.3	$32 \pm 0.20$
Kaiming Initialization	$81.33 \pm 14.10$	$67.35 \pm 26.10$
Gaussian(0, 0.01)	$\textbf{96.18} \pm \textbf{0.23}$	$\textbf{96.22} \pm \textbf{0.22}$

ViT: Architecture: Vit-Small-192 (2.7M parameters), Dataset: CIFAR-10, Flow Hidden Dim: 384, Time Embed Dim: 64

MLP: Architecture: MLP (26.5K parameters), Dataset: MNIST, Flow Hidden Dim: 256, Time Embed Dim: 64 Dropout: 0.1

# G DATASET GENERATION

Table 9 and Table 10 provide the details of the architecture and training hyperparameters used to create the trained neural network datasets that were used to train DeepWeightFlow. The training datasets can be made available on request.

Table 9: Hyperparameters for training the neural networks that were used as the training datasets for Deep-WeightFlow.

Model	Dataset	Params	LR Schedule	Optimizer	LR	Weight Decay	Batch Size	Epochs
MLP	Iris	131	None	Adam	1e-3	0	16	100
MLP	MNIST	26.5K	None	Adam	1e-3	0	64	5
MLP	Fashion	118K	None	AdamW	1e-3	0	128	25
ResNet-18	CIFAR-10	11.2M	Cosine	SGD	0.1	5e-4	128	100
ResNet-20	CIFAR-10	0.27M	None	Adam	1e-3	0	128	5
Vit-Small-192	CIFAR-10	2.8M	Cosine	AdamW	3e-4	0.05	128	300

Table 10: Model architectures for the neural networks used to train DeepWeightFlow. For the MLPs, the first number in the Architecture definition is the input dimension. For the ResNets, "blocks" refer to residual blocks.

Model	Architecture	Parameters	Dataset	Input Dim
MLP	[4, 16, 3]	131	Iris	$4 \times 150$
MLP	[784, 32, 32, 10]	26,506	MNIST	$28 \times 28$
MLP	[784, 128, 128, 10]	117,770	Fashion-MNIST	$28 \times 28$
ResNet-20	$3 \times [3, 3, 3]$ blocks	272,474	CIFAR-10	$32 \times 32 \times 3$
ResNet-18	$4 \times [2, 2, 2, 2]$ blocks	11.17M	CIFAR-10	$32 \times 32 \times 3$
Vit-Small-192	194 embedding dimension, 6 blocks, 3 heads	2.87M	CIFAR-10	$32 \times 32 \times 3$

## H HYPERPARAMETERS OF DEEPWEIGHTFLOW MODELS

In Table 11 we provide the hyperparameters of the DeepWeightFlow models. The FM model architecture varies by the dimensionality of the neural networks weights in the training set and their architecture.

Table 11: DeepWeightFlow Flow Matching training hyperparameters

Parameter	Value	Parameter	Value
Architecture		Training	
Flow Model Hidden Dims	$[d_h, d_h/2, d_h]^a$	Optimizer	AdamW
Time Embedding Dim	4–128 <sup>b</sup>	Learning Rate	$5 \times 10^{-4}$
Activation Function	GELU	Weight Decay	$1 \times 10^{-5}$
Layer Normalization	Yes	AdamW $\beta$	(0.9, 0.95)
Dropout Rate	$0.1 - 0.4^{c}$	Batch Size	4-8 <sup>d</sup>
Flow Matching		Training	
Time Distribution	Uniform	Training Iterations	30,000
Noise Scale $(\sigma)$	0.001	Training Data Size	100 models
Source Distribution	$\mathcal{N}(0,\sigma_s^2 I)^{e}$	LR Scheduler	CosineAnnealing
		$\eta_{ m min}$	$1 \times 10^{-6}$
Generation		Preprocessing	
ODE Solver	Runge-Kutta 4	Weight Matching	Git Re-Basin/TransFusion <sup>f</sup>
Integration Steps	100	BN Recalibration	ResNets only <sup>g</sup>
Generated Samples	100		•

 $<sup>^{\</sup>rm a}$   $d_h \in \{32, 64, 128, 256, 384, 512\}$  depending on architecture complexity

## I COMPUTATIONAL EFFICIENCY: TRAINING AND GENERATION TIME

Table 12: Performance comparison between DeepWeightFlow and RPG (Wang et al., 2025). RPG generates a single neural network per run, while DeepWeightFlow generates 100 neural networks sequentially in a single workflow. For 100 neural networks, we estimate the generation time for RPG based on the time required to generate a single neural network, as presented in Wang et al. (2025).

Model	Method	Hidden Dim	Training Time	Generation Time (100 models)	GPU
	RPG (sequential) <sup>†</sup>	-	-	31 hours	H100
	RPG (partially parallel) <sup>†</sup>	-	-	3 hours	H100
ResNet-18	RPG (fully parallel) <sup>†</sup>	-	-	2.8 hours	H100
(11.7M params)	DeepWeightFlow §	1024	3 min	2.3 min	A100
	DeepWeightFlow + rebasin <sup>§</sup>	1024	$2 \min + 3 \min$	2.3 min	A100
	RPG (flatten) <sup>‡</sup>	-	6.2 hours	16.3 hours	H100
	RPG (by channel) <sup>‡</sup>	-	14.2 hours	16.3 hours	H100
ViT-Tiny	RPG (within layer) <sup>‡</sup>	-	6.2 hours	16.3 hours	H100
(5M params)	RPG (partially parallel) <sup>†</sup>	-	-	1.8 hours	H100
	RPG (fully parallel) <sup>†</sup>	-	-	1.8 hours	H100
	DeepWeightFlow §	256	21 min	3.6 min	A100
Vit-Small-192	DeepWeightFlow §	384	19 min	2.83 min	H100
(2.8M params)	DeepWeightFlow + transfusion <sup>§</sup>	384	13 min + 19 min	2.83 min	H100

<sup>&</sup>lt;sup>†</sup> RPG inference times from Wang et al. (2025) are available only for generating a single model from which we have computed the time necessary to generate 100 models.

DeepWeightFlow demonstrates significant computational advantages over existing parameter generation methods. We compare our approach with RPG (Wang et al., 2025), the current state-of-the-art in recurrent parameter generation, across multiple architectures and configurations.

b Time embedding: 4 for Iris MLP, 64 for ResNet-20/MNIST/Fashion-MNIST/Vit-Small-192, 128 for ResNet-18

<sup>&</sup>lt;sup>c</sup> Dropout: 0.4 for Iris MLP, 0.1 for all other architectures

<sup>&</sup>lt;sup>d</sup> Batch size: 4 for Vit-Small-192 and 8 for all others

 $<sup>^{\</sup>rm e}$   $\sigma_s=0.001$  for Vit-Small-192,  $\sigma_s=0.01$  for all other architectures

f Git Re-Basin for ResNets/MLPs, TransFusion for Vision Transformers

g BatchNorm statistics recalibrated using test data only for ResNet architectures post-generation

<sup>&</sup>lt;sup>h</sup> Generated samples: 25 for Vit-Small-192, 100 for all other architectures

<sup>&</sup>lt;sup>‡</sup> RPG training + sequential inference time from Wang et al. (2025), for generating 100 neural networks estimated from the numbers available for single neural network generation.

<sup>§</sup> DeepWeightFlow performs sequential generation of 100 models.

When incorporating Git-Rebasin (Ainsworth et al., 2023) for weight alignment, the additional computational overhead is minimal:

- ResNet-18: 2 minutes for aligning 100 models
- Vit-Small-192 (Transfusion): 13 minutes for aligning 100 models

The results show that DeepWeightFlow consistently generates high-quality models while having lower training and inference time on similar GPUs.

## J DIVERSITY OF THE GENERATED NEURAL NETWORKS

In Table 13 we provide the numerical estimates of mIoU, the Jensen-Shannon, Wasserstein, and Nearest Neighbors (NN) distances between generated and original neural networks highlighting the diversity of the generated neural networks

Table 13: Comparison of 100 complete MNIST classifying MLP weights generated by DeepWeightFlow with and without Git Re-Basin through maximum Intersection over Union (IoU), Jensen-Shannon, Wasserstein, and Nearest Neighbors (NN) distances. Lower scores show closer relationships. (Org. - original, Gen. - generated)

	Org. to Org.	Org. to Gen.	Gen. to Org.	Gen. to Gen.
DeepWeightFlow w/ Re-Basin				
IoU	-	-	$0.8187 \pm 0.0385$	-
Wasserstein	-	13.4125	21.2867	11.6721
Jensen-Shannon	-	0.7146	0.8326	0.7146
NN	$23.0393 \pm 0.2214$	$9.7232 \pm 10.4398$	$1.7526 \pm 0.1671$	$11.7407 \pm 10.547$
DeepWeightFlow w/o Re-Basin				
IoU	-	-	$0.8256 \pm 0.0748$	-
Wasserstein	-	15.1185	25.6979	17.6939
Jensen-Shannon	-	0.8181	0.8326	0.7293
NN	$27.4895 \pm 0.2007$	$12.3710 \pm 12.441$	$1.7916 \pm 0.3753$	$9.7956 \pm 11.2484$