# Finding Symmetry in Neural Network Parameter Spaces

**Bo Zhao**
University of California, San Diego
bozhao@ucsd.edu

**Nima Dehmamy**
IBM Research
nima.dehmamy@ibm.com

**Robin Walters**
Northeastern University
r.walters@northeastern.edu

**Rose Yu**
University of California, San Diego
roseyu@ucsd.edu

## Abstract

Parameter space symmetries, or loss-invariant transformations, are important for understanding neural networks' loss landscape, training dynamics, and generalization. However, identifying the full set of these symmetries remains a challenge. In this paper, we formalize data-dependent parameter symmetries and derive their infinitesimal form, which enables an automated approach to discover symmetry across different architectures. Our framework systematically uncovers parameter symmetries, including previously unknown ones. We also prove that symmetries in smaller subnetworks can extend to larger networks, allowing the discovery of symmetries in small architectures to generalize to more complex models.

## 1 Introduction

Parameter space symmetry, or loss-invariant transformation of parameters, influences various aspects of deep learning theory. Continuous symmetry connects groups to their orbits, revealing important topological properties such as the dimension [Zhao et al., 2023b] and connectedness [Zhao et al., 2023a] of the minimum. Parameter symmetry also influences training dynamics through the associated conserved quantities of gradient flow [Kunin et al., 2021] and by steering stochastic gradient descent towards certain favored solutions [Ziyin, 2024]. Additionally, symmetry provides a tool to perform optimization within a loss level set, with successful applications in accelerating optimization [Armenta et al., 2023, Zhao et al., 2022] and improving generalization [Zhao et al., 2024]. Other applications of parameter space symmetry include model compression [Ganev et al., 2022, Sourek et al., 2021] and reducing the search space for efficient sampling in Bayesian neural networks [Wiese et al., 2023].

Despite the wide range of applications, our knowledge of parameter space symmetries is limited. In particular, known symmetries often cannot account for all loss-invariant parameter transformation. While several frameworks have been developed to unify known symmetries, whether the symmetries in current literature are complete remains an open question. Due to a lack of systematic approach, current practice typically requires deriving symmetries from scratch for every new architecture, creating barriers for wider application that leverages parameter symmetries. In this paper, we discuss an automated approach to directly learn the symmetry groups and their group actions on the parameter space of neural networks. We show that large networks often have symmetries inherited from its components or subnetworks. This view suggests that searching for symmetries in small networks is an effective approach to identify a significant number of symmetries in modern architectures.

Our main contributions are:

- Formal definitions of data-dependent parameter symmetries and their infinitesimal form.

- An approach to identify symmetries in the parameter space of large networks from known symmetries in smaller subnetworks.
- A framework that discovers symmetry in neural network parameter spaces.

## 2 Data-dependent group action and symmetry

Let $\Theta$ be the space of parameters and $\mathcal{D}$ be the space of data. In this paper, we consider loss functions of the form $L : \Theta \times \mathcal{D} \to \mathbb{R}$, which map parameters and a single data point to a real number. By abuse of notation, we allow $L$ to simultaneously process multiple data points. Specifically, we sometimes define $L : \Theta \times \mathcal{D}^d \to \mathbb{R}^d$ for $d \in \mathbb{Z}^+$ data points.

Let $G$ be a group. Consider a map $a$, which defines a map for every data batch of size $d \in \mathbb{Z}^+$:

$$a : \mathcal{D}^d \to (G \times \Theta \to \Theta)$$
$$X \mapsto (a_X : g, \theta \mapsto \theta'). \tag{1}$$

The map $a$ is a group action on $\Theta$ if it satisfies the following axioms:

identity: $\quad a_X(I, \theta) = \theta, \quad \forall X \in \mathcal{D}^d, \ \forall \theta \in \Theta.$

associative law: $\quad a_X(g_2, a_X(g_1, \theta)) = a_X(g_2 g_1, \theta), \quad \forall g_1, g_2 \in G, \ \forall X \in \mathcal{D}^d, \ \forall \theta \in \Theta.$

A group action $a$ is a parameter space symmetry of $L$ if it additionally satisfies

loss invariance: $\quad L(a_X(g, \theta), X) = L(\theta, X), \quad \forall g \in G, \ \forall X \in \mathcal{D}^d, \ \forall \theta \in \Theta.$

A function $L$ has a $G$-symmetry if there exists a loss-invariant group action $a$. We refer to $G$ as a symmetry group of $L$. Additionally, the action $a$ is termed a data-dependent group action or symmetry if the map (1) has a non-trivial dependency on $X$. That is, $a$ is data-dependent if there exists $X_1, X_2 \in \mathcal{D}^d$, such that $a_{X_1} \neq a_{X_2}$. We derive an infinitesimal version of parameter space symmetries in Appendix B.

## 3 Building symmetries from known ones

One way to identify symmetries in a large network is by examining its components or subnetworks. Despite often having billions of parameters, neural networks typically consist of a limited set of functional families, such as fully connected layers, attention mechanisms, and activation functions. This modular view suggests a mechanism by which symmetries in networks with fewer layers might extend to those in deeper networks. Additionally, within similar types of networks, it may be possible to extrapolate symmetries found in narrower layers to wider ones.
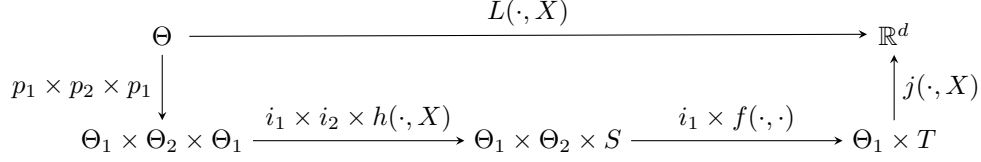
By focusing on symmetries in small architectures and using them to infer symmetries in larger ones, we circumvent the complexity associated with direct handling of high-dimensional parameter spaces. This approach not only simplifies the discovery of symmetries in large-scale networks but also provides a systematic method for using symmetries in smaller subnetworks to understand those in more extensive architectures. Proofs and further discussions can be found in Appendix C.

When a loss function $L$ depends on a subset of parameter exclusively through a subnetwork $f$, any symmetries that preserves $f$ will also preserve the original network $L$:

**Proposition 3.1.** *Let $L : \Theta \times \mathcal{D}^d \to \mathbb{R}^d$ be a function, where the parameter space $\Theta$ is a product space $\Theta = \Theta_1 \times \Theta_2$, with spaces $\Theta_1, \Theta_2$. Suppose there exist functions $h : \Theta_1 \times \mathcal{D}^d \to S$, $f : \Theta_2 \times S \to T$, and $j : (\Theta_1 \times T) \times \mathcal{D}^d \to \mathbb{R}^d$, such that for every $\theta = (\theta_1, \theta_2) \in \Theta$ and $X \in \mathcal{D}^d$, $L(\theta, X) = j\big((\theta_1, f(\theta_2, h(\theta_1, X))), X\big)$. If $a : S \to (G \times \Theta_2 \to \Theta_2)$ is a $G$-symmetry of $f$, then there is an induced $G$-symmetry of $L$, $a' : \mathcal{D}^d \to (G \times \Theta \to \Theta)$, defined by $a'_X(g, (\theta_1, \theta_2)) = \big(\theta_1, a_{h(\theta_1, X)}(g, \theta_2)\big)$.*

The relationship between the functions in the proposition is described by the commutative diagram below, where $p_1 : \Theta \to \Theta_1$, $p_2 : \Theta \to \Theta_2$ are projections onto $\Theta_1$ and $\Theta_2$, $i_1 : \Theta_1 \to \Theta_1$ and $i_2 : \Theta_2 \to \Theta_2$ are identity maps, and $X \in \mathcal{D}^d$ represents a batch of data. When $L$ can be decomposed

in this way, the function $h$ does not depend on $\Theta_2$, and the function $j$ depends on $\Theta_2$ only through the output of $f$. This effectively confines $L$'s dependency on $\Theta_2$ to the transformation defined by $f$, ensuring that any transformation on $\Theta_2$ not altering the output of $f$ will not affect the output of $L$. Consequently, symmetries identified in the smaller network $f$ can be extrapolated to the larger network $L$.

$$
\begin{array}{ccc}
\Theta & \xrightarrow{\;\;\;\;\;\;\;\;\;\;\;\; L(\cdot, X) \;\;\;\;\;\;\;\;\;\;\;\;} & \mathbb{R}^d \\
{\scriptstyle p_1 \times p_2 \times p_1} \downarrow & & \uparrow {\scriptstyle j(\cdot, X)} \\
\Theta_1 \times \Theta_2 \times \Theta_1 \xrightarrow{i_1 \times i_2 \times h(\cdot, X)} \Theta_1 \times \Theta_2 \times S & \xrightarrow{i_1 \times f(\cdot, \cdot)} & \Theta_1 \times T
\end{array}
$$

Proposition 3.1 can be applied to construct symmetries in larger networks from those in smaller ones (Corollary C.2, C.1 in Appendix C). Figure 1 shows the subset of parameters ($\Theta_2$) the symmetry applies to in the corollaries.
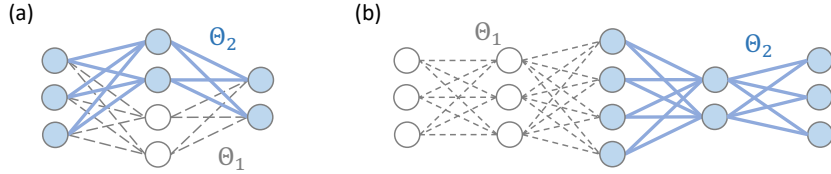


Figure 1: If a network contains substructures with known symmetry, we can infer the same symmetry for the large network. (a) Symmetry from narrower networks. (b) Symmetry from shallower networks.

Note that this approach does not explore the emergence of new, more complex symmetries that may arise as the neural network scale up in size. Notably, there are cases where there exists a $G$ symmetry over its input space, but group actions on individual layers are not loss-invariant (Kvinge et al. [2022]). Nevertheless, studying smaller and simpler networks remains a effective strategy to obtain a significant number of symmetries in larger networks, and is a first step in characterizing the complete set of symmetries in modern architectures.

## 4 Automatic Discovery of Parameter Symmetries

Formulating symmetries in the infinitesimal form makes them easier to learn using an automatic framework, as it defines a set of local conditions for a function to be a symmetry. Using the infinitesimal symmetry derived in Section B.1, we construct an automated framework for discovering parameter space symmetries.

**Enforcing Loss Invariance and Group Axioms.** Given a function $L$, our goal is to find a symmetry $a$ and a set of Lie algebra elements $h$ corresponding to a symmetry group of $L$. We parameterize $a$ using a neural network with learnable parameters, and set $h$ to be learnable as well. We define the following loss terms that quantify the deviation from loss invariance and the group axioms (identity and associativity law):

$$L_{\text{invariance}} = \mathbb{E}_{x,\theta} |D_\theta L|_{\theta,X} \circ D_g a_X|_{I,\theta}(h)| \tag{2}$$
$$L_{\text{id}} = \mathbb{E}_{x,\theta} \|a_x(I, \theta) - \theta\|_2 \tag{3}$$

The two loss terms bias the action towards being loss-invariant and preserving identity. By minimizing $L_{\text{Lie\_deriv}}$, we ensure that the learned symmetry $a$ and the Lie algebra element $h$ satisfy the infinitesimal symmetry condition (Theorem B.1). Minimizing $L_{\text{id}}$ enforces the identity axiom. By focusing on the Lie algebras, we enforce the loss invariance and group structure at the infinitesimal level. This formulation allows us to avoid computing exponential maps.

**Regularizations.** To prevent the group action to be the identity function, we encourage the infinitesimal action to be nonzero. In implementation, we include the following regularization term

to encourage the norm of the infinitesimal action to be around a fixed positive real number $\beta$: $L_{\text{reg\_id}} = \min_{a,h} \mathbb{E}_\theta |\beta - \|D_g a_X|_{I,\theta}(h)\||$.

When learning multiple generators simultaneously, we want them to be orthogonal. Following Yang et al. [2023b], we do this by including the following cosine similarity between each pair of the $k$ generators in the loss function: $L_{\text{reg\_h\_orth}} = \sum_{1 \leq i < j \leq k} \frac{h_i \cdot h_j}{\|h_i\| \|h_j\|}$.

Finally, we encourage sparsity of $h$ for easier interpretation, with $L_{\text{reg\_h\_sparse}} = \sum_{k,j} |h_{kj}|$.

The final training objective is a weighted average of the above loss and regularization terms, with hyperparameters $\gamma_1, ..., \gamma_6 \in \mathbb{R}^+$:

$$\min_{h,a} \gamma_1 L_{\text{invariance}} + \gamma_2 L_{\text{id}} + \gamma_3 L_{\text{reg\_id}} + \gamma_4 L_{\text{reg\_h\_orth}} + \gamma_5 L_{\text{reg\_h\_sparse}}. \quad (4)$$

### 4.1 Learned data-independent symmetries

In the first set of tasks, we see if our method can learn generators for architectures with already known data-independent symmetries. We consider two-layer networks in the form of $L(W_1, W_2, X, Y) = \|W_2 \sigma(W_1 X) - Y\|^2$, where $W_2 \in \mathbb{R}^{m \times h}, W_1 \in \mathbb{R}^{h \times n}$ are parameters, $X \in \mathbb{R}^{n \times k}, Y \in \mathbb{R}^{m \times k}$ are data, and $\sigma$ is a homogeneous activation function.

During training, we train the generators $h$ and the group action $a$ under objective (4). We parametrize $a$ using a 4-layer MLP with hidden dimensions 64, 64, 64. The group aciton $a$ takes a group element, parameter, and data as input and outputs transformed parameters. We use 10000 training samples, each containing a randomly generated set of parameters and data. We set the learning rate as $10^{-3}$ with decay 0.6 every 1000 steps, and the weights for the multi-objective loss as $\gamma_1 = 10, \gamma_2 = \gamma_4 = \gamma_5 = 1$, and $\gamma_6 = 0.1$.



Figure 2: Generator for a two-layer linear MLP with scalar parameters and data.

As a proof of concept, we training a group action and a single generator $h \in \mathbb{R}^{2 \times 2}$ for the two-layer architecture with $m = h = n = k = 1$ and $\sigma$ being the identity function. Figure 2 visualizes the learned generator, which matches the expected generator that generates the rescaling group.

Note that, however, we do not impose constraints on the group action (in particular, not enforcing linear actions). Hence we do not expect the learned generators to look similar to the elements of the Lie algebra infinitesimal generators of the symmetry group in general. For example, the action $a$ can be a composition of two function, the first transforming learned generators to the set of actual generators, and the second performing the group action. We find that our method can learn the generators and group actions for wider two-layer homogeneous architectures as well. More examples of learned generators for larger architectures can be found in Appendix D.

### 4.2 Learned data-dependent Symmetries

As a more practical application of our framework, we attempt to uncover data-dependent symmetries from architectures where no continuous symmetry is known before. We apply our framework to learn generators and loss-invariant group actions for two-layer neural network with sigmoid and tanh activation function, as well as a three-layer neural network with skip connection.

Specifically, we aim to learn symmetries in the two-layer networks defined in the previous section, but replacing $\sigma$ by sigmoid or tanh. Our objective is again to find a set of generators $h$ and a group action $a$ that minimizes (4). We use 10000 training samples, each containing a randomly generated set of parameters and data. We set the learning rate as $10^{-3}$ and the weights for the multi-objective loss as $\gamma_1 = 1, \gamma_2 = \gamma_4 = 10, \gamma_5 = 1$, and $\gamma_6 = 0.1$.

Figure 3 shows the learned generators for data-dependent symmetries in a two-layer sigmoid MLP with parameters dimensions $W_1 \in \mathbb{R}^{3 \times 3}, W_2 \in \mathbb{R}^{3 \times 1}$ and data $X \in \mathbb{R}^{3 \times 1}, Y \in \mathbb{R}^{1 \times 1}$. Figure 6 in the Appendix shows the training curve. Since sigmoid networks have no data-independent continuous symmetry, this set of symmetries are data-dependent, indicating that our method successfully learns data-dependent symmetries for this architecture.
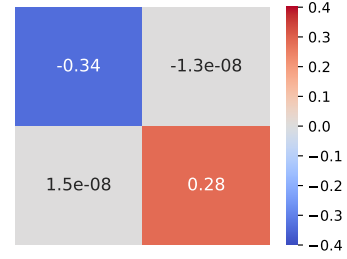
Figure 4 shows the learned generators for data-dependent symmetries in a three-layer tanh MLP with parameters dimensions $W_1 \in \mathbb{R}^{2\times 2}, W_2 \in \mathbb{R}^{2\times 2}, W_3 \in \mathbb{R}^{2\times 1}$ and data $X \in \mathbb{R}^{1\times 2}, Y \in \mathbb{R}^{1\times 1}$. The generators indicate the existence of symmetries that act on non-contiguous layers, which has not been discovered in previous literature.
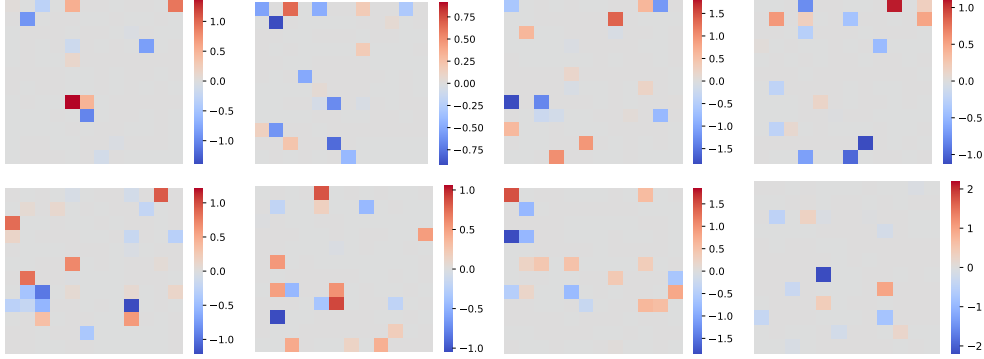


Figure 3: Learned generators for data-dependent symmetries in a two-layer sigmoid MLP with parameters dimensions $W_2 \in \mathbb{R}^{3\times 1}, W_1 \in \mathbb{R}^{3\times 3}$ and data $X \in \mathbb{R}^{1\times 3}, Y \in \mathbb{R}^{1\times 1}$.
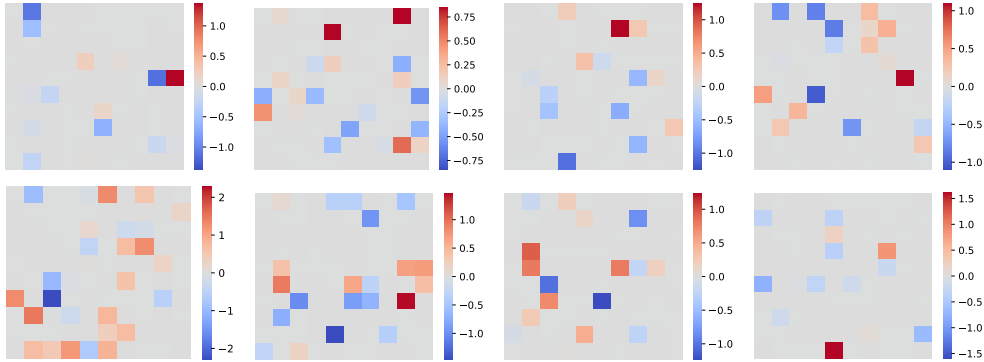


Figure 4: Learned generators for data-dependent symmetries in a three-layer tanh MLP with parameters dimensions $W_1 \in \mathbb{R}^{2\times 2}, W_2 \in \mathbb{R}^{2\times 2}, W_3 \in \mathbb{R}^{2\times 1}$ and data $X \in \mathbb{R}^{1\times 2}, Y \in \mathbb{R}^{1\times 1}$.

## 5  Discussion

While our discovery framework suggests that there are previously unknown data-dependent symmetries in various neural network architectures, the existence and number of symmetries in neural network parameter spaces remain open questions. Whether the number of symmetries is affected by existence of symmetry in data or changes during training are also interesting directions. Future work will examine the structure of learned symmetry, such as the dimension of Lie algebras.

## References

Marco Armenta, Thierry Judge, Nathan Painchaud, Youssef Skandarani, Carl Lemaire, Gabriel Gibeau Sanchez, Philippe Spino, and Pierre-Marc Jodoin. Neural teleportation. *Mathematics*, 11 (2):480, 2023.

Vijay Badrinarayanan, Bamdev Mishra, and Roberto Cipolla. Symmetry-invariant optimization in deep networks. *arXiv preprint arXiv:1511.01754*, 2015.

Gregory Benton, Marc Finzi, Pavel Izmailov, and Andrew G Wilson. Learning invariances in neural networks from training data. *Advances in neural information processing systems*, 33:17605–17616, 2020.

An Mei Chen, Haw-minn Lu, and Robert Hecht-Nielsen. On the geometry of feedforward neural network error surfaces. *Neural computation*, 5(6):910–927, 1993.

Nima Dehmamy, Robin Walters, Yanchen Liu, Dashun Wang, and Rose Yu. Automatic symmetry discovery with lie algebra convolutional network. *Advances in Neural Information Processing Systems*, 34:2503–2515, 2021.

Simon S Du, Wei Hu, and Jason D Lee. Algorithmic regularization in learning deep homogeneous models: Layers are automatically balanced. *Neural Information Processing Systems*, 2018.

Alex Gabel, Victoria Klein, Riccardo Valperga, Jeroen SW Lamb, Kevin Webster, Rick Quax, and Efstratios Gavves. Learning lie group symmetry transformations with neural networks. In *Topological, Algebraic and Geometric Learning Workshops 2023*, pages 50–59. PMLR, 2023.

Iordan Ganev, Twan van Laarhoven, and Robin Walters. Universal approximation and model compression for radial neural networks. *arXiv preprint arXiv:2107.02550v2*, 2022.

Elisenda Grigsby, Kathryn Lindsey, and David Rolnick. Hidden symmetries of relu networks. In *International Conference on Machine Learning*, pages 11734–11760. PMLR, 2023.

Nate Gruver, Marc Anton Finzi, Micah Goldblum, and Andrew Gordon Wilson. The lie derivative for measuring learned equivariance. In *The Eleventh International Conference on Learning Representations*, 2022.

Pavan Karjol, Rohan Kashyap, Aditya Gopalan, and A. P. Prathosh. A unified framework for discovering discrete symmetries. In *Proceedings of The 27th International Conference on Artificial Intelligence and Statistics*, volume 238 of *Proceedings of Machine Learning Research*, pages 793–801. PMLR, 2024.

Sven Krippendorf and Marc Syvaeri. Detecting symmetries with neural networks. *Machine Learning: Science and Technology*, 2(1):015010, 2020.

Daniel Kunin, Javier Sagastuy-Brena, Surya Ganguli, Daniel LK Yamins, and Hidenori Tanaka. Neural mechanics: Symmetry and broken conservation laws in deep learning dynamics. In *International Conference on Learning Representations*, 2021.

Henry Kvinge, Tegan Emerson, Grayson Jorgenson, Scott Vasquez, Tim Doster, and Jesse Lew. In what ways are deep neural networks invariant and how should we measure this? *Advances in Neural Information Processing Systems*, 35:32816–32829, 2022.

Artem Moskalev, Anna Sepliarskaia, Ivan Sosnovik, and Arnold Smeulders. Liegg: Studying learned lie group generators. *Advances in Neural Information Processing Systems*, 35:25212–25223, 2022.

Artem Moskalev, Anna Sepliarskaia, Erik J Bekkers, and Arnold WM Smeulders. On genuine invariance learning without weight-tying. In *Topological, Algebraic and Geometric Learning Workshops 2023*, pages 218–227. PMLR, 2023.

Vasco Portilheiro. Quantifying lie group learning with local symmetry error. In *NeurIPS 2023 Workshop on Symmetry and Geometry in Neural Representations*, 2023.

David W Romero and Suhas Lohit. Learning partial equivariances from data. *Advances in Neural Information Processing Systems*, 35:36466–36478, 2022.

Ben Shaw, Abram Magner, and Kevin R Moon. Symmetry discovery beyond affine transformations. *arXiv preprint arXiv:2406.03619*, 2024.

Sho Sonoda, Hideyuki Ishi, Isao Ishikawa, and Masahiro Ikeda. Joint group invariant functions on data-parameter domain induce universal neural networks. In *NeurIPS 2023 Workshop on Symmetry and Geometry in Neural Representations*, 2023.

Gustav Sourek, Filip Zelezny, and Ondrej Kuzelka. Lossless compression of structured convolutional models via lifting. In *International Conference on Learning Representations*, 2021.

Alonso Urbano and David W Romero. Self-supervised detection of perfect and partial input-dependent symmetries. *arXiv preprint arXiv:2312.12223*, 2023.

Jonas Gregor Wiese, Lisa Wimmer, Theodore Papamarkou, Bernd Bischl, Stephan Günnemann, and David Rügamer. Towards efficient mcmc sampling in bayesian neural networks by exploiting symmetry. In *Joint European Conference on Machine Learning and Knowledge Discovery in Databases (ECML PKDD): Research Track*, pages 459–474, 2023.

Jianke Yang, Nima Dehmamy, Robin Walters, and Rose Yu. Latent space symmetry discovery. *arXiv preprint arXiv:2310.00105*, 2023a.

Jianke Yang, Robin Walters, Nima Dehmamy, and Rose Yu. Generative adversarial symmetry discovery. *International Conference on Machine Learning*, 2023b.

Bo Zhao, Nima Dehmamy, Robin Walters, and Rose Yu. Symmetry teleportation for accelerated optimization. *Advances in Neural Information Processing Systems*, 2022.

Bo Zhao, Nima Dehmamy, Robin Walters, and Rose Yu. Understanding mode connectivity via parameter space symmetry. In *UniReps: the First Workshop on Unifying Representations in Neural Models*, 2023a.

Bo Zhao, Iordan Ganev, Robin Walters, Rose Yu, and Nima Dehmamy. Symmetries, flat minima, and the conserved quantities of gradient flow. *International Conference on Learning Representations*, 2023b.

Bo Zhao, Robert M Gower, Robin Walters, and Rose Yu. Improving convergence and generalization using parameter symmetries. *International Conference on Learning Representations*, 2024.

Allan Zhou, Tom Knowles, and Chelsea Finn. Meta-learning symmetries by reparameterization. *International Conference on Learning Representations*, 2021.

Liu Ziyin. Symmetry leads to structured constraint of learning. In *International Conference on Machine Learning*. PMLR, 2024.

# Appendix

## A   Related Work

**Parameter space symmetry.**   Parameter symmetries are loss-invariant transformations on neural network parameters, often in the form of group actions. Symmetry is present in many neural networks. Known symmetries include invertible linear transformations in linear networks, rescaling in homogeneous networks [Badrinarayanan et al., 2015, Du et al., 2018], radial rescaling in radial neural networks [Ganev et al., 2022], and translation in softmax and scaling in batchnorm functions [Kunin et al., 2021]. In tanh neural networks [Chen et al., 1993], only permutation and sign flip symmetries preserve the loss function. ReLU networks, however, possess symmetries beyond the well-known rescaling [Grigsby et al., 2023]. The existence and number of symmetries in most other architectures remain an open question.

**Data-dependent symmetry.**   While the above symmetries leave the loss unchanged on all data, a relaxed definition, data-dependent symmetry, only requires loss invariance on a subset of data. Zhao et al. [2023b] found examples of such symmetries with nontrivial data dependency, although these symmetries are complicated, limited to minibatches of size one, and difficult generalize across different architectures. This motivates an automated symmetry discovery framework, which, in principle, can find symmetries of arbitrary form in arbitrary architectures. The concept of a symmetry dependent on data has also appeared in adjacent fields. For example, [Moskalev et al., 2023] observe that learned data invariance in neural networks is strongly conditioned on data and breaks under data distribution drift; Sonoda et al. [2023] define a joint group action on data and parameters as part of a new proof of universal approximation theory.

**Discovering and measuring symmetry.**   Various work explores learning continuous symmetries by identifying generators of Lie groups [Krippendorf and Syvaeri, 2020, Moskalev et al., 2022, Dehmamy et al., 2021, Yang et al., 2023b, Gabel et al., 2023], including cases with nonlinear group actions [Yang et al., 2023a, Shaw et al., 2024]. We build on this approach to discover data-dependent group action in high-dimensional parameter spaces. While learning discrete symmetry [Zhou et al., 2021, Karjol et al., 2024] and distributions of symmetry [Benton et al., 2020, Romero and Lohit, 2022, Urbano and Romero, 2023] are also relevant, they are not the primary focus of this paper.

Extracted symmetry is often evaluated locally, by measuring function changes under infinitesimal symmetry transformations [Gruver et al., 2022] or by comparing tangent spaces of orbits under the learned group and the true symmetry group [Portilheiro, 2023]. We adopt the local invariance of loss functions under symmetry transformation, similar to that defined in [Gruver et al., 2022, Moskalev et al., 2022], as the minimization objective in learning data-dependent group actions.

## B   Infinitesimal Symmetry and Examples

### B.1   Infinitesimal Symmetry

We derive an infinitesimal version of parameter space symmetries. For the automatic symmetry discovery framework in Section 4, this definition allows us to learn the group elements and actions without computing matrix exponential, which is expensive, during training.

The following theorem shows that the derivative of the loss function $L$ with respect to the parameters $\theta$ vanishes in the directions generated by the symmetry group's infinitesimal transformations. In other words, the loss function is invariant to small changes along these symmetric directions in parameter space.

**Theorem B.1.** *Let $a : \mathcal{D}^d \to (G \times \Theta \to \Theta)$ be a parameter space symmetry of a loss function $L : \Theta \times \mathcal{D}^d \to \mathbb{R}^d$. Let $D_\theta L|_{\theta,X} : T_\theta\Theta \to \mathbb{R}^d$ be the derivative of $L$ with respect to $\theta$, and $D_g a_X|_{I,\theta} : \mathfrak{g} \to T_\theta\Theta$ be the derivative of $a_X(g, \theta)$ with respect to $g$. Then, for all $\theta \in \Theta$, $X \in \mathcal{D}^d$, and $h \in \mathfrak{g}$,*

$$D_\theta L|_{\theta,X} \circ D_g a_X|_{I,\theta} \circ h = 0. \tag{5}$$

*Proof sketch.* Consider a smooth curve $\gamma(t) = a_X(\exp(ht), \theta)$ in $\Theta$, where $h \in \mathfrak{g}$ and $t \in \mathbb{R}$. Then, since $L$ is invariant under $a$, $L(\gamma(t), X) = L(\theta, X), \forall t \in \mathbb{R}$. The result follows from differentiating both sides with respect to $t$ at $t = 0$ and applying chain rules. $\qquad\square$

*Proof.* Since $a$ is a symmetry of $L$, we have

$$L(a_X(g, \theta), X) = L(\theta, X), \quad \forall g \in G, \quad \forall \theta \in \Theta, \quad \forall X \in \mathcal{D}^d.$$

Consider a smooth curve $\gamma(t) = a_X(\exp(ht), \theta)$ in $\Theta$, where $h \in \mathfrak{g}$ and $t \in \mathbb{R}$. Then, since $L$ is invariant under $a$,

$$L(\gamma(t), X) = L(\theta, X), \quad \forall t \in \mathbb{R}.$$

Differentiating both sides with respect to $t$ at $t = 0$, we get

$$\frac{d}{dt} L(\gamma(t), X)\bigg|_{t=0} = 0.$$

Applying the chain rule,

$$\frac{d}{dt} L(\gamma(t), X)\bigg|_{t=0} = D_\theta L|_{\theta, X} \left( \frac{d\gamma(t)}{dt}\bigg|_{t=0} \right).$$

Now, compute $\frac{d\gamma(t)}{dt}\big|_{t=0}$ using the chain rule:

$$\frac{d\gamma(t)}{dt}\bigg|_{t=0} = \frac{d}{dt} a_X(\exp(ht), \theta)\bigg|_{t=0} = D_g a_X|_{I, \theta} \left( \frac{d}{dt} \exp(ht)\bigg|_{t=0} \right).$$

Since $\exp$ is the exponential map from $\mathfrak{gl}(n)$ to $GL(n)$, and $h \in \mathfrak{gl}(n)$, we have

$$\frac{d}{dt} \exp(ht)\bigg|_{t=0} = h.$$

Therefore,

$$\frac{d\gamma(t)}{dt}\bigg|_{t=0} = D_g a_X|_{I, \theta}(h).$$

Putting it all together,

$$D_\theta L|_{\theta, X} \left( D_g a_X|_{I, \theta}(h) \right) = 0.$$

$\qquad\square$

Equation 5 states that the gradient of the loss function $L$ with respect to the parameters $\theta$ is orthogonal to the directions in parameter space generated by the infinitesimal symmetry transformations $D_g a_X|_{I, \theta}(h)$. This orthogonality implies that moving along these symmetric directions does not change the loss to first order, reflecting the invariance of $L$ under the group action.

Assuming that $\Theta = \mathbb{R}^n$, then for a single data point ($d = 1$), we can write (5) in coordinates as

$$D_\theta L|_{\theta, X} \left( D_g a_X|_{I, \theta}(h) \right) = \sum_{i=1}^{dim(\Theta)} \sum_{k=1}^{dim(\mathfrak{g})} \frac{\partial L}{\partial \theta_i} \left( D_g a_X|_{I, \theta} \right)_{ik} h_k = 0. \tag{6}$$

When $\Theta = \mathbb{R}^n$ and $G$ is a subgroup of $GL(n)$ with a linear, data-independent symmetry $a_x(g, \theta) = g\theta$ for all $x \in X$, (6) reduces to the equation in Theorem 3.1 in [Moskalev et al., 2022]. With $(D_g a)_{ijk} = \frac{\partial a_i}{\partial g_{jk}} = \delta_{ij}\theta_k$, we have

$$\frac{dL(\exp(h \cdot t) \cdot \theta)}{dt}\bigg|_{t=0} = \sum_{i=1}^{n} \sum_{j=1}^{n} \sum_{k=1}^{n} \frac{\partial L}{\partial \theta_i} \left( D_g a|_{I, \theta} \right)_{ijk} h_{jk} = \sum_{i=1}^{n} \sum_{k=1}^{n} \frac{\partial L}{\partial \theta_i} \theta_k h_{ik}. \tag{7}$$

Our symmetry acts on parameters instead of data, but otherwise this matches Theorem 3.1 in [Moskalev et al., 2022].

### B.2 Alternative Option for Discovery Objectives

A more straightforward training objective exponentiates the Lie algebra to obtain group elements, before enforcing loss invariance and group axioms:

$$\min_{h,a} L_{\text{invariance\_int}} + L_{\text{id\_int}} + L_{\text{assoc\_int}}$$

with

$$L_{\text{invariance\_int}} = \mathbb{E}_{x,\theta,t}\|L\left(a_x(exp(ht),\theta),x\right) - L(\theta,x)\|$$
$$L_{\text{id\_int}} = \mathbb{E}_{x,\theta}\|a_x(I,\theta) - \theta\|$$
$$L_{\text{assoc\_int}} = \sum_{h_1,h_2\in\mathfrak{g}} \mathbb{E}_{x,\theta}\left\|a_{\exp(h_1)X}(\exp(h_2), a_X(\exp(h_1),\theta)) - a_X(\exp(h_2)\exp(h_1),\theta)\right\|.$$

Similarly to the infinitesimal version, this objective also directly enforces the necessary group structures. We adopt the infinitesimal formulation to avoid the computational overhead of evaluating exponential maps.

## C  Building symmetries from known ones

This section contains the proofs for results in Section 3.

**Proposition 3.1.** *Let $L : \Theta \times \mathcal{D}^d \to \mathbb{R}^d$ be a function, where the parameter space $\Theta$ is a product space $\Theta = \Theta_1 \times \Theta_2$, with spaces $\Theta_1, \Theta_2$. Suppose there exist functions $h : \Theta_1 \times \mathcal{D}^d \to S$, $f : \Theta_2 \times S \to T$, and $j : (\Theta_1 \times T) \times \mathcal{D}^d \to \mathbb{R}^d$, such that for every $\theta = (\theta_1,\theta_2) \in \Theta$ and $X \in \mathcal{D}^d$, $L(\theta,X) = j\big((\theta_1, f\big(\theta_2, h(\theta_1,X)\big)),X\big)$. If $a : S \to (G \times \Theta_2 \to \Theta_2)$ is a $G$-symmetry of $f$, then there is an induced $G$-symmetry of $L$, $a' : \mathcal{D}^d \to (G \times \Theta \to \Theta)$, defined by $a'_X(g, (\theta_1,\theta_2)) = \big(\theta_1, a_{h(\theta_1,X)}(g,\theta_2)\big)$.*

*Proof.* We need to show that $a'$ satisfies the identity and associative law of a group action and preserves $L$.

Since $a$ is a group action on $\Theta_2$, it satisfies the identity axiom $a_{h(\theta_1,X)}(I,\theta_2) = \theta_2$. Applying this in the definition of $a'$, we get $a'_X(I, (\theta_1,\theta_2)) = (\theta_1, a_{h(\theta_1,X)}(I,\theta_2)) = (\theta_1,\theta_2)$.

Since $a$ is a group action on $\Theta_2$, it satisfies the associative law $a_{h(\theta_1,X)}(g_2 g_1, \theta_2) = a_{h(\theta_1,X)}(g_2, a_{h(\theta_1,X)}(g_1, \theta_2))$, for all $g_1, g_2 \in G$. It follows that $a'$ also satisfies the associative law: $a'_X(g_2 g_1, (\theta_1,\theta_2)) = (\theta_1, a_{h(\theta_1,X)}(g_2 g_1, \theta_2)) = (\theta_1, a_{h(\theta_1,X)}(g_2, a_{h(\theta_1,X)}(g_1, \theta_2))) = a'_X(g_2, a'_X(g_1, (\theta_1,\theta_2)))$

Finally, since $a$ is a symmetry of $f$, we have $f(a_{h(\theta_1,X)}(g,\theta_2), h(\theta_1,X)) = f(\theta_2, h(\theta_1,X))$, for all $g \in G$. It follows that $a'$ preserves the value of $L$: $L(a'_X(g,\theta),X) = j\big((\theta_1, f\big(a_{h(\theta_1,X)}(g,\theta_2), h(\theta_1,X)\big)),X\big) = j\big((\theta_1, f\big(\theta_2, h(\theta_1,X)\big)),X\big) = L(\theta,X)$. $\square$

The first corollary describes how symmetries identified in narrower networks also apply to wider networks. A function $\sigma : \mathbb{R}^{h\times k} \to \mathbb{R}^{h\times k}$ is row-wise if, for any matrix $A \in \mathbb{R}^{h\times k}$ with rows $\{a_i \in \mathbb{R}^k\}_{i=1}^h$, the output matrix $\sigma(A)$ has rows $\{\sigma_{row}(a_i) \in \mathbb{R}^k\}_{i=1}^h$, where $\sigma_{row} : \mathbb{R}^k \to \mathbb{R}^k$ applies independently on each row of $A$. Element-wise function is a special case of row-wise functions. For fully connected networks with row-wise activation functions, identifying a symmetry in one architecture suggests that the same symmetry will apply to wider versions of that architecture.

**Corollary C.1.** *Consider a network parameter space $\Theta(m,h,n) = \mathbb{R}^{m\times h} \times \mathbb{R}^{h\times n}$ and data space $\mathcal{D}(n,k) = \mathbb{R}^{n\times k}$. Let $\sigma : \mathbb{R}^{h\times k} \to \mathbb{R}^{h\times k}$ be a row-wise function. Consider a function $L_{mnhk} : \Theta(m,h,n) \times \mathcal{D}(n,k) \to \mathbb{R}^{m\times k}$, defined as $L_{mnhk}((U,V),X) = U\sigma(VX)$ for $U \in \mathbb{R}^{m\times h}$, $V \in \mathbb{R}^{h\times n}$, and $X \in \mathbb{R}^{n\times k}$. If there is a $G$-symmetry of $L_{mnhk}$, then there is a $G$-symmetry of $L_{mnh'k}$ with any $h' > h$.*

*Proof.* The function $L_{mnh'k}$ can be decomposed into

$$U(\sigma(VX))_{ik} = U_{ij}\sigma(VX)_{jk}$$

$$= \sum_{j=1}^{h}\sum_{l=1}^{n} U_{ij}\sigma(V_{jl}X_{lk})$$

$$= \sum_{j=1}^{h}\sum_{l=1}^{n} U_{ij}\sigma(V_{jl}X_{lk}) + \sum_{j=h+1}^{h'}\sum_{l=1}^{n} U_{ij}\sigma(V_{jl}X_{lk}) \qquad (8)$$

Note that for all $i, k$, the first term depends only on the first $h$ columns of $U$ and first $h$ rows of $V$, and the second terms depends only on the rest of the columns and rows of $U$ and $V$. Denoting the first $h$ columns of $U$ as $U_{1:h}$, the rest of the columns of $U$ as $U_{h+1:h'}$, the first $h$ rows of $V$ as $V_{1:h}$, and the rest of the rows of $V$ as $V_{h+1:h'}$, we have

$$L_{mnh'k}((U,V),X) = L_{mnhk}((U_{1:h},V_{1:h}),X) + L_{mn(h'-h)k}((U_{h+1:h'},V_{h+1:h'}),X). \qquad (9)$$

Let $\Theta_1 = \mathbb{R}^{m\times h}\times\mathbb{R}^{h\times n}$ and $\Theta_2 = \mathbb{R}^{m\times(h'-h)}\times\mathbb{R}^{(h'-h)\times n}$. Then $\Theta(m,h',n) = \Theta_1\times\Theta_2$. Let $S = (\mathbb{R}^{m\times k}\times\mathcal{D}^d)$ and $T = \mathbb{R}^{m\times k}\times\mathbb{R}^{m\times k}$. Define the following three functions

$$h : \Theta_1\times\mathcal{D}^d \to (\mathbb{R}^{m\times k}\times\mathcal{D}^d)$$
$$f : \Theta_2\times(\mathbb{R}^{m\times k}\times\mathcal{D}^d) \to \mathbb{R}^{m\times k}\times\mathbb{R}^{m\times k}$$
$$j : (\Theta_1\times(\mathbb{R}^{m\times k}\times\mathbb{R}^{m\times k}))\times\mathcal{D}^d \to \mathbb{R}^{m\times k} \qquad (10)$$

by

$$h((U_{1:h},V_{1:h}),X) = (L_{mnhk}((U_{1:h},V_{1:h}),X),X)$$
$$f((U_{h+1:h'},V_{h+1:h'}),(Y,X)) = \big(L_{mn(h'-h)k}((U_{h+1:h'},V_{h+1:h'}),X),Y\big)$$
$$j\big(((U_{1:h},V_{1:h}),(Y',Y)),X\big) = Y'+Y. \qquad (11)$$

Then $L_{mnh'k}(\theta,X) = j\big((\theta_1,f(\theta_2,h(\theta_1,X))),X\big)$ for all $\theta = (\theta_1,\theta_2)\in\Theta$ and $X\in\mathcal{D}^d$. Since $L_{mnhk}$ has a symmetry, $f$ has the same symmetry. By Proposition 3.1, $L_{mnh'k}$ also has the same symmetry. $\qquad\square$

The next corollary shows that symmetries of a subset of layers are also symmetries in the entire network. Both corollary can be proved by partitioning the dimensions the parameter space and defining corresponding functions that compose $L$, before applying Proposition 3.1. Figure 1 shows the subset of parameters ($\Theta_2$) the symmetry applies to in the corollaries. These are the subnetworks where symmetries are assume to be known, which the larger network inherits.

**Corollary C.2.** *Let $\Theta = \Theta_1\times...\times\Theta_l$ be a parameter space. Consider a list of spaces $\Phi_0 = \mathcal{D}^d$, $\Phi_l = \mathbb{R}^d$, and $\Phi_1, ..., \Phi_{l-1}$. Let $L : \Theta\times\mathcal{D}^d\to\mathbb{R}^d$ be a function defined recursively by $\{L_i\}_{i=1}^{l}$ with $L_i : \Theta_i\times\Phi_{i-1}\to\Phi_i$, such that $L = \phi_l$ where $\phi_i = L_i(\theta_i,\phi_{i-1})\in\Phi_i$ and $\phi_0 = X$. If for some $1\le i\le l$, $L_i$ has a $G$-symmetry, then $L$ has a $G$-symmetry.*

*Proof.* Define functions

$$h : (\Theta_1\times...\times\Theta_{i-1}\times\Theta_{i+1}\times...\times\Theta_l)\times\mathcal{D}^d \to \Phi_{i-1}$$
$$f : \Theta_i\times\Phi_{i-1} \to \Phi_i$$
$$j : (\Theta_1\times...\times\Theta_{i-1}\times\Theta_{i+1}\times...\times\Theta_l)\times\Phi_i\times\mathcal{D}^d \to \mathbb{R}^d \qquad (12)$$

by

$$h((\theta_1,...,\theta_{i-1},\theta_{i+1},...,\theta_l),X) = L_{i-1}(\theta_{i-1},X), \quad \text{computed using } (\theta_1,...,\theta_{i-1})$$
$$f(\theta_i,\phi_{i-1}) = L_i(\theta_i,\phi_{i-1})$$
$$j((\theta_1,...,\theta_{i-1},\theta_{i+1},...,\theta_l),\phi_i,X) = L_l(\theta_l,X), \quad \text{computed using } (\theta_l,...,\theta_{i+1}) \text{ and } \phi_i. \qquad (13)$$

Then $L((\theta_1,...,\theta_l),X) = j\big((\theta_1,...,\theta_{i-1},\theta_{i+1},...,\theta_l),f(\theta_i,h((\theta_1,...,\theta_{i-1},\theta_{i+1},...,\theta_l),X)),X\big)$ for all $\theta = (\theta_1,\theta_2)\in\Theta$ and $X\in\mathcal{D}^d$. By Proposition 3.1, if $f = L_i$ has a $G$-symmetry, $L$ also has a $G$-symmetry. $\qquad\square$

In addition to obtaining symmetries from those in smaller networks, we can also get symmetries for a loss function over data batches with a certain size, if we know there is a symmetry for this function over larger data batches. Concretely, if there exists a group action that preserves loss for all data batches of size $d \in \mathbb{Z}^+$, then that group action preserves loss for all data batches of size $d' < d$.

**Proposition C.3.** *Let $L_d : \Theta \times \mathcal{D}^d \to \mathbb{R}^d$ be a function that is applied pointwise on each of $d$ data points in a data batch. If $L_d$ admits a $G$-symmetry, then $L_{d'}$ admits a $G$-symmetry for all $d' < d$.*

*Proof.* Suppose that $L_d$ has a $G$-symmetry. Let $a : \mathcal{D}^d \to (G \times \Theta \to \Theta), X_d \mapsto (a_{X_d} : g, \theta \mapsto \theta')$ be the corresponding group action. Define $a' : \mathcal{D}^{d'} \to (G \times \Theta \to \Theta)$ by $X_{d'} \mapsto (a_{t(X_{d'})} : g, \theta \mapsto \theta')$, where $t : \mathcal{D}^{d'} \to \mathcal{D}^d$ appends $d - d'$ random data points to its input. Clearly, $a'$ satisfies the identity and associate axiom and preserves loss. Therefore, $a'$ is a $G$-symmetry of $L_{d'}$. ☐

# D   Additional Experiment Details
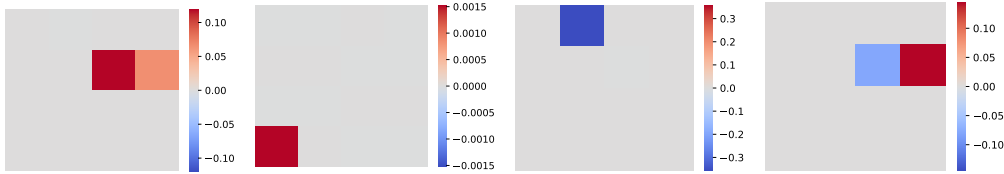


Figure 5: Learned generators for a two-layer linear MLP with parameters dimensions $W_2 \in \mathbb{R}^{1 \times 2}, W_1 \in \mathbb{R}^{2 \times 1}$ and data $X, Y \in \mathbb{R}$.
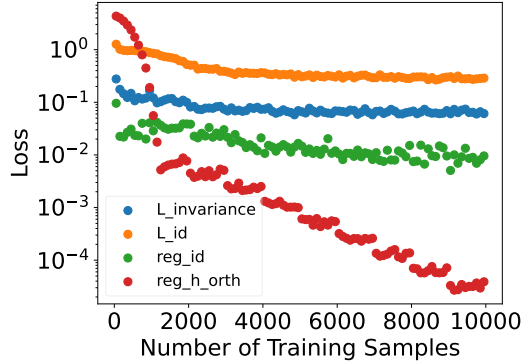


Figure 6: Training curve for learning data-dependent symmetry in a two-layer sigmoid MLP with parameters dimensions $W_2 \in \mathbb{R}^{3 \times 1}, W_1 \in \mathbb{R}^{3 \times 3}$ and data $X \in \mathbb{R}^{1 \times 3}, Y \in \mathbb{R}^{1 \times 1}$.