
Automatic Symmetry Discovery with Lie Algebra Convolutional Network

Nima Dehmamy
Northwestern University
nimadt@bu.edu

Robin Walters
Northeastern University
rwalters@northeastern.edu

Yanchen Liu
Northeastern University

Dashun Wang
Northwestern University
dashun.wang@kellogg.northwestern.edu

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

Abstract

Existing equivariant neural networks require prior knowledge of the symmetry group and discretization for continuous groups. We propose to work with Lie algebras (infinitesimal generators) instead of Lie groups. Our model, the Lie algebra convolutional network (L-conv) can automatically discover symmetries and does not require discretization of the group. We show that L-conv can serve as a building block to construct *any* group equivariant feedforward architecture. Both CNNs and Graph Convolutional Networks can be expressed as L-conv with appropriate groups. We discover direct connections between L-conv and physics: (1) group invariant loss generalizes field theory (2) Euler-Lagrange equation measures the robustness, and (3) equivariance leads to conservation laws and Noether current. These connections open up new avenues for designing more general equivariant networks and applying them to important problems in physical sciences.

1 Introduction

Incorporating symmetries into a deep learning architecture can reduce sample complexity, improve generalization, while significantly decreasing the number of model parameters (Cohen et al., 2019b; Cohen & Welling, 2016b; Ravanbakhsh et al., 2017; Ravanbakhsh, 2020; Wang et al., 2020). For instance, Convolutional Neural Networks (CNN) (LeCun et al., 1989, 1998) implement translation symmetry through weight sharing. General principles for constructing symmetry-aware group equivariant neural networks were introduced in Cohen & Welling (2016b), Kondor & Trivedi (2018), and Cohen et al. (2019b).

However, most work on equivariant networks requires knowing the symmetry group *a priori*. A different equivariant model needs to be re-designed for each symmetry group. In practice, we may not have a good inductive bias and such knowledge of the symmetries may not be available. Constructing and selecting the equivariant network with the appropriate symmetry group becomes quite tedious. Furthermore, many existing works are limited to *finite groups* such as permutations Hartford et al. (2018); Ravanbakhsh et al. (2017); Zaheer et al. (2017), 90 degree rotations Cohen et al. (2018) or dihedral groups D_N $E(2)$ Weiler & Cesa (2019).

For a continuous group, existing approaches either discretize the group Weiler et al. (2018a,b); Cohen & Welling (2016a), or use a truncated sum over irreducible representations (irreps) Weiler & Cesa (2019); Weiler et al. (2018a) via spherical harmonics in Worrall et al. (2017) or more general Clebsch-Gordon coefficients Kondor et al. (2018); Bogatskiy et al. (2020). These approaches are prone to approximation error. Recently, Finzi et al. (2020) propose to approximate the integral over the Lie group by Monte Carlo sampling. This approach requires implementing the matrix exponential and obtaining a local neighborhood for each point. Both parametrizing Lie groups for sampling and finding irreps are computationally expensive. Finzi et al. (2021) provide a general algorithm for constructing equivariant multi-layer perceptrons (MLP), but require explicit knowledge of the group to encode its irreps, and solving a set of constraints.

We provide a novel framework for designing equivariant neural networks. We leverage the fact that Lie groups can be constructed from a set of infinitesimal generators, called Lie algebras. A Lie algebra has a finite basis, assuming the group is finite-dimensional. Working with the Lie algebra basis allows us to encode an infinite group without discretizing or summing over irreps. Additionally, all Lie algebras have the same general structure and hence can be implemented the same way. We propose Lie Algebra Convolutional Network (**L-conv**), a novel architecture that can automatically discover symmetries from data. Our main contributions can be summarized as follows:

We propose the Lie algebra convolutional network (**L-conv**), a building block for constructing group equivariant neural networks.

We prove that multi-layer L-conv can approximate group convolutional layers, including CNNs, and find graph convolutional networks to be a special case of L-conv.

We can learn the Lie algebra basis in L-conv, enabling automatic symmetry discovery.

L-conv also reveals interesting connections between physics and learning: equivariant loss generalizes important Lagrangians in field theory; robustness and equivariance can be expressed as Euler-Lagrange equations and Noether currents.

Learning symmetries from data has been studied in limited settings for commutative Lie groups as in Cohen & Welling (2014), 2D rotations and translations in Rao & Ruderman (1999), Sohl-Dickstein et al. (2010) or permutations (Anselmi et al., 2019). (Zhou et al., 2020) propose a general method for symmetry discovery. Yet, their weight-sharing scheme and the symmetry generators are very different from ours. Our approach use much fewer parameters and has a direct interpretation using Lie algebras (SI B.3). Benton et al. (2020) propose Augerino to learn a distribution over data augmentations. It also involves Lie algebras, but is restricted to a subgroup of 2D affine transformations and requires matrix logarithm and sampling (SI B.3). In contrast, our approach is simpler and more general.

2 Background

We review the core concepts L-conv builds upon: equivariance, group convolution and Lie algebras.

Notations. Unless explicitly stated, a in A^a is an index, not an exponent. We use the Einstein summation $A^a B_{ab} = \sum_a A^a B_{ab} = [AB]_b$, where a repeated upper and lower index are summed.

Equivariance. Let S be a topological space on which a Lie group G (continuous group) acts from the left, meaning for all $\mathbf{x} \in S$ and $g \in G$, $g\mathbf{x} \in S$. We refer to S as the base space. Let F , the “feature space”, be the vector space $F = \mathbb{R}^m$. Each data point is a feature map $f : S \rightarrow F$. The action of G on the input of f induces an action on feature maps. For “scalar” features, for $u \in G$, the transformed features $u \cdot f$ are given by

$$u \cdot f(\mathbf{x}) = f(u^{-1} \mathbf{x}). \tag{1}$$

Denote the space of all functions from S to F by F^S , so that $f \in F^S$. Let F be a mapping to a new feature space $F^0 = \mathbb{R}^{m^0}$, meaning $F : F^S \rightarrow F^0$. We say F is *equivariant* under G if G acts on F^0 and for $u \in G$, we have

$$u \cdot (F(f)) = F(u \cdot f). \tag{2}$$

Group Convolution. Kondor & Trivedi (2018) showed that F is a linear equivariant map if and only if it performs a group convolution (G-conv). To define G-conv, we first lift \mathbf{x} to elements in G (Kondor & Trivedi, 2018). Specifically, we pick an origin $\mathbf{x}_0 \in S$ and replace each point $\mathbf{x} = g\mathbf{x}_0$ by

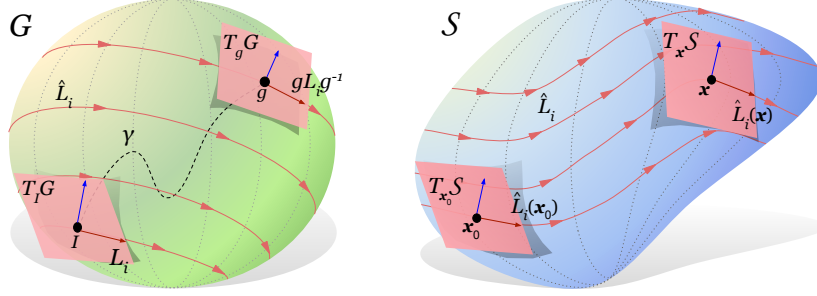


Figure 1: **Lie group and Lie algebra:** Illustration of the group manifold of a Lie group G (left). The Lie algebra $\mathfrak{g} = T_I G$ is the tangent space at the identity I . L_i are a basis for $T_I G$. If G is connected, $\forall g \in G$ there exist paths like γ from I to g and g can be written as a path-ordered integral $g = P \exp[\int_{\gamma} dt^i L_i]$. **Base space** Right is a schematic of the base space S as a manifold. The lift $\mathbf{x} = g\mathbf{x}_0$ takes $\mathbf{x} \in S$ to $g \in G$, and maps the tangent spaces $T_{\mathbf{x}} S \rightarrow T_g G$. Each Lie algebra basis $L_i \in \mathfrak{g} = T_I G$ generates a vector field \hat{L}_i on the tangent bundle TG via the pushforward $\hat{L}_i(g) = gL_i g^{-1}$. Via the lift, L_i also generates a vector field $\hat{L}_i = \hat{L}_i^\alpha(\mathbf{x})\partial_\alpha = [gL_i \mathbf{x}_0]^\alpha \partial_\alpha$.

g . We will often drop \mathbf{x}_0 for brevity and write $f(g) = f(g\mathbf{x}_0)$. Let $\kappa : G \rightarrow \mathbb{R}^{m^0} \rightarrow \mathbb{R}^m$ be a linear transformation from F to F^0 . G-conv is defined as

$$[\kappa \star f](g) = \int_G \kappa(g^{-1}v) f(v) dv = \int_G \kappa(v) f(gv) dv, \quad (3)$$

We denote the Haar measure on G as $dv = d\mu(v)$ for brevity.

Equivariance of G-conv. G-conv in equation 3 is equivariant (Kondor & Trivedi, 2018). By definition, for $w \in G$ we have

$$\begin{aligned} [\kappa \star w \star f](g) &= \int_G \kappa(v) w \star f(gv) dv = \int_G \kappa(v) f(w^{-1}gv) dv \\ &= [\kappa \star f](w^{-1}g) = w \star [\kappa \star f](g) \end{aligned} \quad (4)$$

Existing works on equivariance networks implement \int_G by discretizing the group or summing over irreps. We take a different approach and use the infinitesimal generators of the group. While a Lie group G is infinite, usually it can be generated using a small number of infinitesimal generator, comprising its ‘‘Lie algebra’’. We use the Lie algebra to introduce a building block to approximate G-conv. Figure 1 visualizes a Lie group, Lie algebra and the concept we discuss below.

Lie algebra. Let G be a Lie group, which includes common continuous groups. Group elements $u \in G$ infinitesimally close to the identity element I can be written as $u = I + \epsilon^i L_i$ (note Einstein summation), where $L_i \in \mathfrak{g}$ with the Lie algebra $\mathfrak{g} = T_I G$ is the tangent space of G at the identity element. The Lie algebra has the property that it is closed under a Lie bracket $[\cdot, \cdot] : \mathfrak{g} \rightarrow \mathfrak{g}$

$$[L_i, L_j] = c_{ij}^k L_k, \quad (5)$$

which is skew-symmetric and satisfies the Jacobi identity. Here the coefficients $c_{ij}^k \in \mathbb{R}$ or \mathbb{C} are called the structure constants of the Lie algebra. For matrix representations of \mathfrak{g} , $[L_i, L_j] = L_i L_j - L_j L_i$ is the commutator. The L_i are called the infinitesimal generators of the Lie group.

Exponential map. If the manifold of G is connected¹, an exponential map $\exp : \mathfrak{g} \rightarrow G$ can be defined such that $g = \exp[t^i L_i] \in G$. For matrix groups, if G is connected and compact, the matrix exponential is such a map and it is surjective. For most other groups (except $GL_d(\mathbb{C})$ and nilpotent groups) it is not surjective. Nevertheless, for any connected group every $g \in G$ can be written as a product $g = \int_a \exp[t_a^i L_i]$ (Hall, 2015). Making t_a^i infinitesimal steps $dt^i(s)$ tangent to a path γ from I to g on G yields the surjective path-ordered exponential in physics, denoted as $g = P \exp[\int_{\gamma} dt^i L_i]$ (SI A, and see Time-ordering in Weinberg (1995, p143)).

¹When G has multiple connected components, these results hold for the component containing I , and generalize easily for multi-component groups such as $Z_k \otimes G$ (Finzi et al., 2021).

Pushforward. $L_i \in T_I G$ can be pushed forward to $\hat{L}_i(g) = gL_i g^{-1} \in T_g G$ to form a basis for $T_g G$, satisfying the same Lie algebra $[\hat{L}_i(g), \hat{L}_j(g)] = c_{ij}^k \hat{L}_k(g)$. The manifold of G together with the set of all $T_g G$ attached to each g forms the tangent bundle TG , a type of fiber bundle (Lee et al., 2009). \hat{L}_i is a vector field on TG . The lift maps \hat{L}_i to an equivalent vector field on TS , which we will also denote by \hat{L}_i . Figure 1 illustrates the flow of these vector fields on TG and TS .

3 Lie Algebra Convolutional Network

We can use the Lie algebra basis $L_i \in \mathfrak{g}$ to construct the Lie group G with the exponential map. Similarly, we show that Lie algebras can also serve as building blocks to construct G-conv layers. We propose the Lie algebra convolutional network (L-conv). The key idea is to approximate the kernel $\kappa(u)$ using localized kernels which can be constructed using the Lie algebra (Fig. 2). This is possible because the exponential map is a generalization of a Taylor expansion. We show that a G-conv whose kernel is concentrated near the identity can be expanded in the Lie algebra.

Let $\delta_\eta(u) \in \mathbb{R}$ denote a normalized *localized kernel*, meaning $\int_G \delta_\eta(g) dg = 1$, and with support on a small neighborhood of size η near the identity I (i.e., $\delta_\eta(I + \epsilon^i L_i) \neq 0$ if $k\epsilon k^2 > \eta^2$). Let $\kappa_0(u) = W^0 \delta_\eta(u)$, where $W^0 \in \mathbb{R}^{m^0} \times \mathbb{R}^m$ are constants. The localized kernels κ_0 can be used to approximate G-conv. We first derive the expression for a G-conv whose kernel is κ_0 .

Linear expansion of G-conv with localized kernel. We can expand a G-conv whose kernel is $\kappa_0(u) = W^0 \delta_\eta(u)$ in the Lie algebra of G to linear order. With $v_\epsilon = I + \epsilon^i L_i$, we have (see SI A)

$$\begin{aligned} Q[f](g) &= [\kappa_0 \star f](g) = \int_G dv \kappa_0(v) f(gv) = \int_{k\epsilon k < \eta} dv_\epsilon \kappa_0(v_\epsilon) f(gv_\epsilon) \\ &= W^0 \int d\epsilon \delta_\eta(v_\epsilon) f(g) + \epsilon^i g L_i \frac{d}{dg} f(g) + O(\epsilon^2) \\ &= W^0 \left[I + \bar{\epsilon}^i g L_i \frac{d}{dg} \right] f(g) + O(\eta^2) \end{aligned} \quad (6)$$

where $W^0 \in \mathbb{R}^{m^0} \times \mathbb{R}^m$, and using $\int_G \delta_\eta(g) dg = 1$ and $\int d\epsilon \delta_\eta(v_\epsilon) = 1$ and $\epsilon^i \in \mathbb{R}^m \times \mathbb{R}^m$, we defined

$$\bar{\epsilon}^i = \int d\epsilon \delta_\eta(v_\epsilon) \epsilon^i \in \mathbb{R}^m \times \mathbb{R}^m. \quad (7)$$

Note that because $k\epsilon k < \eta$ we also have $k\bar{\epsilon} k < \eta$ (SI A, equation 24). Here $d\epsilon$ is the integration measure on the Lie algebra $\mathfrak{g} = T_I G$ induced by the Haar measure dv_ϵ on G .

Interpreting the derivatives. In a matrix representation of G , we have $gL_i \frac{df}{dg} = [gL_i]_\alpha^\beta \frac{df}{dg} = \text{Tr} [gL_i]^T \frac{df}{dg}$. This can be written in terms of partial derivatives $\partial_\alpha f(\mathbf{x}) = \partial f / \partial \mathbf{x}^\alpha$ as follows.

Using $\mathbf{x}^\rho = g_\sigma^\rho \mathbf{x}_0^\sigma$, we have $\frac{df(g\mathbf{x}_0)}{dg} = \mathbf{x}_0^\beta \partial_\alpha f(\mathbf{x})$, and so

$$\hat{L}_i f(\mathbf{x}) = gL_i \frac{df}{dg} = [gL_i]_\beta^\alpha \mathbf{x}_0^\beta \partial_\alpha f(\mathbf{x}) = [gL_i \mathbf{x}_0]^\alpha f(\mathbf{x}) \quad (8)$$

Hence, for each L_i , the pushforward $gL_i g^{-1}$ generates a flow on S through the vector field \hat{L}_i $gL_i \frac{d}{dg} = [gL_i g^{-1} \mathbf{x}]^\alpha \partial_\alpha$ (Fig. 1).

Lie algebra convolutional (L-conv) layer. Equation 6 states that for a kernel localized near the identity, the effect of the kernel can be summarized in W^0 and $\bar{\epsilon}^i \hat{L}_i$. Note that we do not need to perform the integral over G explicitly anymore. Instead of working with a kernel κ_0 , we only need to specify W^0 and $\bar{\epsilon}^i$. Hence, in general, we define the Lie algebra convolution (L-conv) as

$$\begin{aligned} Q[f](\mathbf{x}) &= W^0 \left[I + \bar{\epsilon}^i \hat{L}_i \right] f(\mathbf{x}) \\ &= W^0 \left[I + \bar{\epsilon}^i [gL_i \mathbf{x}_0]^\alpha \partial_\alpha \right] f(\mathbf{x}) \end{aligned} \quad (9)$$

Being an expansion of G-conv, L-conv inherits the equivariance of G-conv, as we show next.

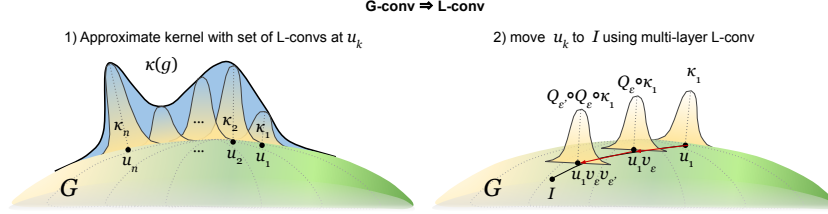


Figure 2: Sketch of the procedure for approximating G-conv using L-conv. First, the kernel is written as the sum of a number of localized kernels κ_k with support around u_k (left). Each of the κ_k is then moved toward identity by composing multiple L-conv layers $Q_{\epsilon^0} \dots Q_{\epsilon} \dots \kappa_k$ (right).

Proposition 1 (Equivariance of L-conv). *With assumptions above, L-conv is equivariant under G .*

Proof: First, note that the components of \hat{L}_i transform as $[\hat{L}_i(v\mathbf{x})]^\alpha = [vgL_i\mathbf{x}_0]^\alpha = v_\beta^\alpha \hat{L}_i(\mathbf{x})^\beta$, while the partial transforms as $\partial/\partial[v\mathbf{x}]^\alpha = [v^{-1}]^\gamma_\alpha \partial_\gamma$. As a result in $\hat{L}_i = [gL_i\mathbf{x}_0]^\alpha \partial_\alpha$ all factors of v cancel, meaning for $v \in G$, $\hat{L}_i(v\mathbf{x}) = \hat{L}_i(\mathbf{x})$. This is because of the fact that $\hat{L}_i \in TS$ is a vector field (i.e. 1-tensor) and, thus, invariant under change of basis. Plugging into equation 9, for $w \in G$

$$\begin{aligned} w \cdot Q[f](\mathbf{x}) &= Q[f](w^{-1}\mathbf{x}) = W^0 \cdot I + \bar{\epsilon}^i \hat{L}_i(w^{-1}\mathbf{x}) \cdot f(w^{-1}\mathbf{x}) \\ &= W^0 \cdot I + \bar{\epsilon}^i \hat{L}_i(g) \cdot f(w^{-1}\mathbf{x}) = W^0 \cdot I + \bar{\epsilon}^i \hat{L}_i(g) \cdot w \cdot f(\mathbf{x}) = Q[w \cdot f](\mathbf{x}) \end{aligned} \quad (10)$$

which proves L-conv is equivariant.

Examples. Using equation 8 we can calculate L-conv for specific groups (details in SI A.2). For translations $G = T_n = (\mathbb{R}^n, +)$, we find the generators become simple partial derivatives $\hat{L}_i = \partial_i$ (SI A.2.2), yielding $f(\mathbf{x}) + \epsilon^\alpha \partial_\alpha f(\mathbf{x})$. For 2D rotations (SI A.2.1) the generator $\hat{L} = (x\partial_y - y\partial_x) = \partial_\theta$, which is the angular momentum operator about the z-axis in quantum mechanics and field theories. For rotations with scaling, $G = SO(2) \times \mathbb{R}^+$, we have two L_i , one $\hat{L}_\theta = \partial_\theta$ from $so(2)$ and a scaling with $L_r = I$, yielding $\hat{L}_r = x\partial_x + y\partial_y = r\partial_r$. Next, we discuss the form of L-conv on discrete data.

3.1 Approximating G-conv using L-conv

L-conv can be used as a basic building block to construct G-conv with more general kernels. Figure 2 sketches the argument described here (see also SI A.1).

Theorem 1 (G-conv from L-convs). *G-conv equation 3 can be approximated using L-conv layers.*

Proof: The procedure involves two steps, as illustrated in Fig. 2: 1) approximate the kernel using localized kernels as the δ_η in L-conv; 2) move the kernels towards identity using multiple L-conv layers. The following lemma outline the details.

Lemma 1 (Approximating the kernel). *Let the kernel $\kappa : G \rightarrow \mathbb{R}$ with $\int_G \kappa(g) k^2 dg < 1$ be continuously differentiable with $k d\kappa(g)/dg k^2 < \xi^2$, and with compact support over $G_0 \subset G$. Let $\kappa_k(g) = c_k \delta_\eta(u_k^{-1}g)$ be a set of N kernels with support on an η neighborhood of $u_k \in G$. Then there exist $c_k \in \mathbb{R}$ and $u_k \in G$ such that $\tilde{\kappa} = \sum_{k=1}^N \kappa_k$ approximates κ , meaning $\int_G \kappa(g) \tilde{\kappa}(g) k^2 dg < \zeta^2$ for arbitrary small $\zeta \in \mathbb{R}_+$.*

Proof: See SI A.1 for details. The intuition is similar to the universal approximation theorem for neural networks (Hornik et al., 1989; Cybenko, 1989), only generalized to a group manifold instead of \mathbb{R} . Let B_0 be the set of $v_\epsilon = I + \epsilon^i L_i \in G$, with $k\epsilon k^2 < \eta^2$. Choose a set of $u_k \in G$ such that the neighborhoods $B_k = u_k B_0 \subset G$ cover the support G_0 of κ . The bound $k d\kappa(g)/dg k^2 < \xi^2$ means that on small enough neighborhoods $B_k \subset G$, for any two $u, v \in B_k$ we have $k\kappa(u) - \kappa(v)k^2 < \eta^2 \xi^2$, where jG_0j is the volume of the support of κ . Hence, for $g \in B_k$, $\kappa(g)$ can be approximated with $\kappa_k(g) = \kappa(u_k) \delta_\eta(u_k^{-1}g)$, with normalized localized kernels $\delta_\eta(g)$, and any element $u_k \in B_k$. We show that the approximation error of using $\tilde{\kappa} = \sum_{k=1}^N \kappa_k$ to approximate κ is bounded by $\int_G \kappa(g) \tilde{\kappa}(g) k^2 dg < jG_0j \eta^2 \xi^2$. Any desired error bound ζ can then be attained by choosing small enough η for neighborhood sizes.

Thus, we can approximate a large class of kernels as $\kappa(g) = \prod_k \kappa_k(g)$ where the local kernels $\kappa_k(g) = c_k \delta_\eta(u_k^{-1}g)$ have support only on an η neighborhood of $u_k \in G$. Here $c_k \in \mathbb{R}^{m^0}$, \mathbb{R}^m are constants and $\delta_\eta(u)$ is as in equation 6. Using this, G-conv equation 3 becomes

$$[\kappa \star f](g) = \prod_k c_k \int \delta_\eta(u_k^{-1}v) f(gv) = \prod_k c_k [\delta_\eta \star f](gu_k). \quad (11)$$

The kernels κ_k are localized around u_k , whereas in L-conv the kernel is around identity. We can compose L-conv layers to move κ_k from u_k to identity.

Lemma 2 (Moving kernels to identity). *κ_k can be moved near identity using a multilayer L-conv.*

Proof: In equation 11, write $u_k = v_\epsilon u_k^0$, with $v_\epsilon = I + \epsilon^i L_i \in \mathfrak{g}$. Using the definition equation 9 an L-conv layer $Q_\epsilon = I - \epsilon^i \hat{L}_i$ performs a first order Taylor expansion (SI A.1) and so $Q_\epsilon[\delta_\eta](u_k^0^{-1}v) = \delta_\eta(u_k^0^{-1}v) + O(\epsilon^2)$. Thus, applying one L-conv layer moves the localized kernel along v_ϵ on G . Writing u_k as the product of a set of small group elements $u_k = \prod_{a=1}^p v_a$, with $v_a = I + \epsilon_a^i L_i \in \mathfrak{g}$. Defining L-conv layers $Q_a = I - \epsilon_a^i \hat{L}_i$, we can write

$$\kappa_k(g) = c_k Q_p \dots Q_1 \delta_\eta(g) \quad (12)$$

meaning κ_k localized around u_k can be written as a p layer L-conv acting on a kernel $\delta_\eta(g)$, localized around the identity of the group. With $\epsilon_a k < \eta$, the error in u_k is $O(\eta^{p+1})$.

Thus, we conclude that any G-conv equation 3 can be approximated by multilayer L-conv. Furthermore, for compact G , using the theorem in Kondor & Trivedi (2018), we can show that any equivariant feedforward neural network can be approximated using multilayer L-conv with nonlinearities.

Equivariance of nonlinearity. Pointwise nonlinearities give equivariant maps between scalar feature maps. To see this, let $\sigma : \mathbb{R} \rightarrow \mathbb{R}$. We extend $\sigma : F \rightarrow F$ by applying σ component-wise. Let $f : S \rightarrow F$ be a scalar feature map (i.e., $g \cdot f(\mathbf{x}) = f(g^{-1}\mathbf{x})$). Then

$$g \cdot (\sigma \circ f)(\mathbf{x}) = \sigma \circ f(g^{-1}\mathbf{x}) = \sigma \circ (g \cdot f)(\mathbf{x}).$$

Since the composition of equivariant maps is equivariant, given equivariant linear mapping $Q : F^S \rightarrow F^{OS}$ (i.e. $g \cdot Q[f] = Q[g \cdot f]$), the layer $f \mapsto \sigma \circ Q[f]$ is equivariant. Hence we have the corollary:

Corollary 1. *Assume G is compact and acts on S transitively. Then any equivariant feedforward neural network (FNN) can be approximated using multilayer L-conv with point-wise nonlinearities.*

Proof: A FNN is defined as $\sigma_p \circ F_p \circ [\sigma_1 \circ F_1[f]](\mathbf{x})$ where F_k are linear and σ_k are point-wise nonlinearities. By Theorem 1 of Kondor & Trivedi (2018), any linear layer in the equivariant FNN is a G-conv, which by Theorem 1 can be approximated by multilayer L-conv. Therefore, multilayer L-conv with nonlinearity can approximate any equivariant FNN.

Finally, to our knowledge it is not known whether *every equivariant function* can be approximated by equivariant FNN for a Lie group G . Hence, the corollary above is *not* a universal approximation theorem for equivariant scalar functions in terms of L-conv. However, it does show that multilayer L-conv is equally expressive as other equivariant networks. Next, we discuss implementation details.

4 Discretized space and implementation: the tensor notation

In many datasets, such as images, $f(\mathbf{x})$ is not given as continuous function, but rather as a discrete array, with $S = \{\mathbf{x}_0, \dots, \mathbf{x}_{d-1}\}$ containing d points. Each \mathbf{x}_μ represents a coordinate in higher dimensional space, e.g. on a 10×10 image, \mathbf{x}_0 is $(x, y) = (0, 0)$ point and \mathbf{x}_{99} is $(x, y) = (9, 9)$.

Feature maps and group action In the tensor notation, we encode $\mathbf{x}_\mu \in S$ as the canonical basis (one-hot) vectors in $\mathbf{x}_\mu \in \mathbb{R}^d$ with $[\mathbf{x}_\mu]_\nu = \delta_{\mu\nu}$ (Kronecker delta), e.g. $\mathbf{x}_0 = (1, 0, \dots, 0)$. The features become $\mathbf{f} \in F = \mathbb{R}^d \times \mathbb{R}^m$, meaning $d \times m$ tensors, with $f(\mathbf{x}_\mu) = \mathbf{x}_\mu^T \mathbf{f} = \mathbf{f}_\mu$. Although S is discrete, the group acting on F can be continuous (e.g. image rotations). Any $G \subset \text{GL}_d(\mathbb{R})$ of the general linear group (invertible $d \times d$ matrices) acts on $\mathbf{x}_\mu \in \mathbb{R}^d$ and $\mathbf{f} \in F$. We define $f(g \cdot \mathbf{x}_\mu) = \mathbf{x}_\mu^T g^T \mathbf{f}$, $\delta g \in G$, so that for $w \in G$ we have

$$w \cdot f(\mathbf{x}_\mu) = f(w^{-1} \cdot \mathbf{x}_\mu) = \mathbf{x}_\mu^T w^{-1T} \mathbf{f} = [w^{-1} \cdot \mathbf{x}_\mu]^T \mathbf{f} \quad (13)$$

Dropping the position \mathbf{x}_μ , the transformed features are matrix product $w \cdot \mathbf{f} = w^{-1T} \mathbf{f}$. We can write G-conv in this notation (SI B). Similarly, we can rewrite L-conv equation 6 in the tensor notation. Defining $v_\epsilon = I + \bar{\epsilon}^i L_i$

$$\begin{aligned} Q[\mathbf{f}](g) &= W^0 \mathbf{f} \cdot g \cdot I + \bar{\epsilon}^i L_i = \mathbf{x}_0^T \cdot I + \bar{\epsilon}^i L_i \cdot g^T \mathbf{f} W^{0T} \\ &= \mathbf{x} + \bar{\epsilon}^i [g L_i \mathbf{x}_0]^T \mathbf{f} W^{0T}. \end{aligned} \quad (14)$$

Here, $\hat{L}_i = g L_i \mathbf{x}_0$ is exactly the matrix analogue of pushforward vector field \hat{L}_i in equation 8. The equivariance of L-conv in tensor notation is again evident from the $g^T \mathbf{f}$, resulting in

$$Q[w \cdot \mathbf{f}](g) = \mathbf{x}_0^T v_\epsilon^T g^T w^{-1T} \mathbf{f} W^{0T} = Q[\mathbf{f}](w^{-1} g) = w \cdot Q[\mathbf{f}](g) \quad (15)$$

Tensor L-conv layer implementation The discrete space L-conv equation 14 can be rewritten using the global Lie algebra basis \hat{L}_i

$$Q[\mathbf{f}] = \mathbf{f} + \hat{L}_i \mathbf{f} \bar{\epsilon}^i W^{0T}, \quad Q[\mathbf{f}]_\mu^a = \mathbf{f}_\mu^b [W^{0T}]_b^a + [\hat{L}_i]_\mu^\nu \mathbf{f}_\nu^c W^i{}_c^a \quad (16)$$

Where $W^i = W^0 \bar{\epsilon}^i$, $W^0 \in \mathbb{R}^{m_{in} \times m_{out}}$ and $\bar{\epsilon}^i \in \mathbb{R}^{m_{in} \times m_{in}}$ are trainable weights. The \hat{L}_i can be either inserted as inductive bias or they can be learned to discover symmetries.

To implement L-conv, note that the formula of equation 16 is quite similar to a Graph Convolutional Network (GCN) (Kipf & Welling, 2016). For each i , the shared convolutional weights are $\bar{\epsilon}^i W^{0T}$ and the aggregation function of the GCN, a function of the graph adjacency matrix, is \hat{L}_i in L-conv. Thus, L-conv can be implemented as GCN modules for each \hat{L}_i , plus a residual connection for the $\mathbf{f} W^{0T}$ term.

Figure 3 shows the schematic of the L-conv layer. In a naive implementation, \hat{L}_i can be general $d \times d$ matrices. However, being vector fields generated by the Lie algebra, \hat{L}_i has a more constrained structure which allows them to be encoded and learned using much fewer parameters than a $d \times d$ matrix. Specifically, encoding the topology of S as a graph (see SI B.1), the incidence matrix replaces partial derivatives (Schaub et al., 2020) in equation 8 and the L_i become weighting of the edges. This weighting is similar to Gauge Equivariant Mesh (GEM) CNN (Cohen et al., 2019a). Indeed, in L-conv the lift $\mathbf{x}_\mu = g_\mu \mathbf{x}_0$ fixes the gauge by mapping neighbors of \mathbf{x}_0 to neighbors of \mathbf{x}_μ . Changing how the discrete S samples an underlying continuous space will change g_μ and hence the gauge.

Choosing the number of L_i . Beside the width of W^0 and $\bar{\epsilon}^i$, the number n_L of L_i is a hyperparameter in L-conv. For instance, if S is a discretization of n dimensional space the symmetry group is likely $G = \text{GL}_n(\mathbb{R}) \cap T_n$, with $n_L = O(n^2)$. Note that n_L is independent of the size d of the discretized space (e.g. number of pixels) and generally $n^2 \gg d$. Choosing n_L larger than the true number of L_i only results in an over-complete basis and shouldn't be a problem. We conducted small controlled experiments to verify how multilayer L-conv approximates G-conv (SI C).

Learning symmetries using L-conv. Rao & Ruderman (1999) introduced a basic version of L-conv and showed that it can learn 1D translation and 2D rotation. We conducted experiments to learn large rotation angle between two images (SI C), shown in Fig. 4. Left shows the architecture for learning the rotation angles between a pair of 7 × 7 random images \mathbf{f} and $R(\theta) \mathbf{f}$ with $\theta \in [0, \pi/3)$. Second left is the learned $L \in \text{SO}(2)$ using 3 recursive layer L-conv. Middle is the L learned using L-conv with fixed small rotation angle $\theta = \pi/10$ (SI C.2) and right is the exact solution $R = (Y X^T)(X^T X)^{-1}$. While the middle L is less noisy, it does not capture weights beyond first neighbors of each pixel. (also see SI C for a discussion on symmetry discovery literature.)

L-conv can potentially replace other equivariant layers in a neural network. We conducted limited experiments for this on small image datasets (SI D). L-conv allows one to look for potential symmetries in data which may have been scrambled or harbors hidden symmetries.

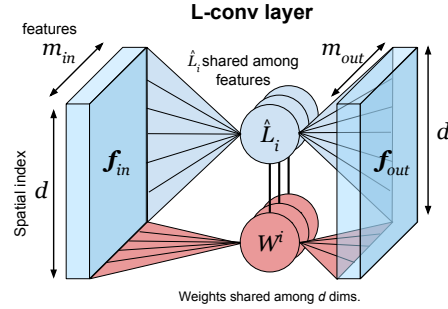


Figure 3: L-conv layer architecture. L_i only act on the d flattened spatial dimensions, and W^i only act on the m_{in} input features and returns m_{out} output features. For each i , L-conv is analogous to a GCN with d nodes and m_{in} features.

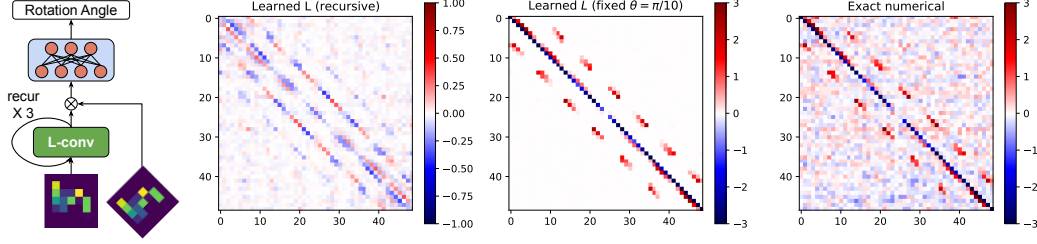


Figure 4: **Learning the infinitesimal generator of $SO(2)$** Left shows the architecture for learning rotation angles between pairs of images (SI C.3). Next to it is the L learned using recursive L-conv in this experiment. Middle L is learned using a fixed small rotation angle $\theta = \pi/10$, and right shows L found using the numeric solution from the data.

5 Relation to other architectures

CNN. This is a special case of expressing G-conv as L-conv when the group is continuous 1D translations. The arguments here generalize trivially to higher dimensions. Rao & Ruderman (1999, sec. 4) used the Shannon-Whittaker Interpolation (Whittaker, 1915) to define continuous translation on periodic 1D arrays as $\mathbf{f}_\rho^\theta = g(z)_\rho^\nu \mathbf{f}_\nu$. Here $g(z)_\rho^\nu = \frac{1}{d} \sum_{p=-d/2}^{d/2} \cos \frac{2\pi p}{d} (z + \rho - \nu)$ approximates the shift operator for continuous z . These $g(z)$ form a 1D translation group G as $g(w)g(z) = g(w+z)$ with $g(0)_\rho^\nu = \delta_\rho^\nu$. For any $z = \mu \in \mathbb{Z}$, $g_\mu = g(z = \mu)$ are circulant matrices that shift by μ as $[g_\mu]_\nu^\rho = \delta_\nu^{\rho - \mu}$. Thus, a 1D CNN with kernel size k can be written using g_μ as

$$F(\mathbf{f})_\nu^a = \sigma \underset{\mu=0}{\star} \mathbf{f}_\nu^c [W^\mu]_\mu^a + b^a = \sigma \underset{\mu=0}{\star} [g_\mu \mathbf{f}]_\nu^c [W^\mu]_\mu^a + b^a \quad (17)$$

where W, b are the filter weights and biases. g_μ can be approximated using the Lie algebra and written as multi-layer L-conv as in sec. 3.1. Using $g(0)_\rho^\nu = \delta(\rho - \nu)$, the single Lie algebra basis $[\hat{L}]_0 = \partial_z g(z)|_{z=0}$, acts as $\hat{L}f(z) = \partial_z f(z)$ (because $\partial_z \delta(z - \nu) f(z) = \partial_\nu f(\nu)$). Its components are $\hat{L}_\rho^\nu = L(\rho - \nu) = \sum_{p=-d/2}^{d/2} \frac{2\pi p}{d} \sin \frac{2\pi p}{d} (\rho - \nu)$, which are also circulant due to the $(\rho - \nu)$ dependence. Hence, $[\hat{L}\mathbf{f}]_\rho = \sum_\nu L(\rho - \nu) \mathbf{f}_\nu = [L \star \mathbf{f}]_\rho$ is a convolution. Rao & Ruderman (1999) already showed that this \hat{L} can reproduce finite discrete shifts g_μ used in CNN. They used a primitive version of L-conv with $g_\mu = (I + \epsilon \hat{L})^N$. Thus, L-conv can approximate 1D CNN. This result generalizes easily to higher dimensions.

Graph Convolutional Network (GCN). Let \mathbf{A} be the adjacency matrix of a graph. In equation 16 if $\hat{L} = h(\mathbf{A})$, such as $\hat{L}_i = \mathbf{D}^{-1/2} \mathbf{A} \mathbf{D}^{-1/2}$, we obtain a GCN (Kipf & Welling, 2016) ($\mathbf{D}_{\mu\nu} = \delta_{\mu\nu} \sum_\rho \mathbf{A}_{\mu\rho}$ being the degree matrix). So in the special case where all neighbors of each node $\langle \mu \rangle$ have the same edge weight, meaning $[\hat{L}_i]_\mu^\nu = [\hat{L}_i]_\mu^\rho, \forall \nu, \rho \in \langle \mu \rangle$, equation 6 is uniformly aggregating over neighbors and L-conv reduces to a GCN. Note that this similarity is not just superficial. In GCN $h(\mathbf{A}) = \hat{L}$ is in fact a Lie algebra basis. When $\hat{L} = h(\mathbf{A})$, the vector field is the flow of isotropic diffusion $d\mathbf{f}/dt = h(\mathbf{A})\mathbf{f}$ from each node to its neighbors. This vector field defines one parameter Lie group with elements $g(t) = \exp[h(\mathbf{A})t]$. Hence, L-conv for flow groups with a single generator are GCN. These flow groups include Hamiltonian flows and other linear dynamical systems. The main difference between L-conv and GCN is that L-conv can assign a different weight to each neighbor of the same node, similar to GEM-CNN (Cohen et al., 2019a) with a fixed gauge set by g_μ . Next, we discuss the mathematical properties of the loss functions for L-conv.

6 Group invariant loss

Loss functions of equivariant networks are rarely discussed. Yet, recent work by Kunin et al. (2020) showed the existence of symmetry directions in the loss landscape. To understand how the symmetry generators in L-conv manifest themselves in the loss landscape, we work out the explicit example of a mean square error (MSE) loss. Because G is the symmetry group, f and $g \star f$ should result in the

same optimal parameters. Hence, the minima of the loss function need to be *group invariant*. One way to satisfy this is for the loss itself to be group invariant, which can be constructed by integrating over G (global pooling (Bronstein et al., 2021)). A function $I = \int_G dg F(g)$ is G -invariant (SI A.3). We can also change the integration to $\int_S d^n x$ by change of variable $dg/d\mathbf{x}$ (see SI A.3 for discussion on stabilizers).

MSE loss and Field Theory. The MSE is given by $I = \int_n \int_G dg k Q[f_n](g) k^2$, where f_n are data samples and $Q[f]$ is L-conv or another G -equivariant function. In supervised learning the input is a pair f_n, y_n . G can also act on the labels y_n . We assume that y_n are either also scalar features $y_n : S \rightarrow \mathbb{R}^{m_y}$ with a group action $g y_n(\mathbf{x}) = y_n(g^{-1}\mathbf{x})$ (e.g. f_n and y_n are both images), or that y_n are categorical. In the latter case $g y_n = y_n$ because the only representations of a continuous G on a discrete set are constant. We can concatenate the inputs to $\phi_n = [f_n/y_n]$ with a well-defined G action $g \phi_n = [g f_n/g y_n]$. The collection of combined inputs $\Phi = (\phi_1, \dots, \phi_N)^T$ is an $(m + m_y) \times N$ matrix. Using equations 6 and 8, the MSE loss with parameters $W = \int W^0, \bar{\epsilon} g$ becomes (SI A.3.1)

$$\begin{aligned} I[\Phi; W] &= \int_G dg L[\Phi; W] = \int_G dg W^0 I + \bar{\epsilon}^i [\hat{L}_i]^\alpha \partial_\alpha \Phi(g)^2 \\ &= \int_S \frac{d^n x}{\frac{\partial x}{\partial g}} \Phi^T \mathbf{m}_2 \Phi + \partial_\alpha \Phi^T \mathbf{h}^{\alpha\beta} \partial_\beta \Phi + [\hat{L}_i]^\alpha \partial_\alpha \Phi^T \mathbf{v}^i \Phi \end{aligned} \quad (18)$$

Equation 18 generalizes the free field theories in physics (Polyakov, 2018). Here $\frac{\partial x}{\partial g}$ is the determinant of the Jacobian, $W^i = W^0 \bar{\epsilon}^i$ and

$$\mathbf{m}_2 = W^{0T} W^0, \quad \mathbf{h}^{\alpha\beta}(\mathbf{x}) = \bar{\epsilon}^{iT} \mathbf{m}_2 \bar{\epsilon}^j [\hat{L}_i]^\alpha [\hat{L}_j]^\beta, \quad \mathbf{v}^i = \mathbf{m}_2 \bar{\epsilon}^i. \quad (19)$$

Note that \mathbf{h} has feature space indices via $[\bar{\epsilon}^{iT} \mathbf{m}_2 \bar{\epsilon}^j]_{ab}$, with index symmetry $\mathbf{h}_{ab}^{\alpha\beta} = \mathbf{h}_{ba}^{\beta\alpha}$. When $F = \mathbb{R}$ (i.e. f is a 1D scalar), $\mathbf{h}^{\alpha\beta}$ becomes a Riemannian metric for S . In general \mathbf{h} combines a 2-tensor $\mathbf{h}_{ab} = \mathbf{h}_{ab}^{\alpha\beta} \partial_\alpha \partial_\beta$ $\int TS \rightarrow TS$ with an inner product $h^T \mathbf{h}^{\alpha\beta} f$ on the feature space F .

In field theory, the motivation is to preserve spatial symmetries for the metric \mathbf{h} . In equation 18, \mathbf{h} transforms equivariantly as a 2-tensor $v \mathbf{h}^{\alpha\beta} = [v^{-1}]^\alpha_\rho [v^{-1}]^\beta_\gamma \mathbf{h}^{\rho\gamma}(\mathbf{x})$ for $v \in G$ (SI A.3). The last term in equation 18 vanishes for many groups (SI A.3) and it is also absent in physics.

Robustness and Euler-Lagrange Equation. Equivariant neural networks are more robust. To check this, we can quantify how the network would perform for an input $\phi^\theta = \phi + \delta\phi$ which adds a small random perturbation $\delta\phi$ to a data point ϕ . Robustness to such perturbation would mean that, for optimal parameters W , the loss function would not change, i.e. $I[\phi^\theta; W] = I[\phi; W]$, requiring I to be minimized around real data points ϕ .

This can be cast as a variational equation $\delta I[\phi; W] = 0$, which yield the familiar Euler-Lagrange (EL) equation (SI A.4). Therefore, for an equivariant network to be robust, i.e. $\delta I[\phi; W]/\delta\phi = 0$, we would require the data points ϕ to satisfy the EL equations for optimal parameters W :

$$\text{Robustness to random noise () EL: } \frac{\partial L}{\partial \phi^b} - \partial_\alpha \frac{\partial L}{\partial (\partial_\alpha \phi^b)} = 0 \quad (20)$$

where the partial derivative terms appear because of the L-conv layer.

Equivariance and Conservation laws. Conserved currents, via Noether's theorem provide a way to find hidden symmetries (see also Kunin et al. (2020)). The idea is that the equivariance condition equation 2 can be written for the integrand of the loss, $L[\phi, W]$. If we write the equivariance equation for infinitesimal v_ϵ , we obtain a vector field which is divergence free. Since G is the symmetry of the system, transforming an input $\phi \rightarrow w \phi$ by $w \in G$ the integrand should change equivariantly, meaning $L[w \phi] = w L[\phi]$. When robustness error is minimized as in equation 20, an infinitesimal $w = I + \eta^i \hat{L}_i$, with $\delta\phi = \epsilon^i \hat{L}_i \phi$, results in a conserved current (SI A.4)

$$\text{Noether current: } J^\alpha = \frac{\partial L}{\partial (\partial_\alpha \phi^b)} \delta \phi^b - \frac{\partial L}{\partial \mathbf{x}^\alpha} \delta \mathbf{x}^\alpha, \quad \delta I[\phi; W] = 0 \quad \Rightarrow \quad \partial_\alpha J^\alpha = 0 \quad (21)$$

The above equation shows that for equivariant networks with a given symmetry, the deviation in data along the symmetry direction (\hat{L}_i) yields a divergence free current J^α , known as Noether current.

It also provides an alternative means to discover symmetry generators L_i by minimizing $k\partial_\alpha J^\alpha k$. Note that this Noether current is the “stress-energy” tensor, associated with space (or space-time) variations $\delta\mathbf{x}$ (Landau, 2013) (SI A.5). We can potentially design more general equivariant networks leading to other Noether currents.

7 Conclusion and Discussions

We propose the Lie algebra convolutional neural network (L-conv), an infinitesimal version of G-conv. L-conv layers do not require encoding irreps or discretizing the group, and can be combined to approximate *any* feedforward equivariant networks on compact groups. Additionally, L-conv’s universal and simple structure allows us to discover symmetries from data. It is easy to implement, with a formula similar to GCN. We validated that L-conv can learn the correct Lie algebra basis in a synthetic experiment.

We discover several intriguing connections between L-conv and physics. Our derivation shows that equivariant neural networks based on L-conv lead to Noether’s theorem and conservation laws. Conversely, we can also optimize Noether current to discover symmetries. Furthermore, the current equivariance formulation only pertains to “spatial symmetries” (i.e. G acts on S). In physics, more general “internal symmetries” are quite common (e.g. particle physics). We can potentially design more general equivariant networks with L-conv encoding such symmetries.

Our method also shed lights on scientific machine learning, especially for physical sciences. Physicists generally use simple polynomial forms for the Lagrangian, or the loss function. These “perturbative” Lagrangian lead to divergences in quantum field theory. However, it is believed the true Lagrangian is more complicated. Hence, more expressive L-conv based models can potentially provide more advanced ansatzes for solving scientific problems.

Acknowledgments and Disclosure of Funding

R. Walters is supported by a Postdoctoral Fellowship from the Roux Institute and NSF grants #2107256 and #2134178. This work was supported in part by the U. S. Army Research Office under Grant W911NF-20-1-0334, DOE ASCR 2493 and NSF Grant #2134274. N. Dehmamy and D. Wang were supported by the Air Force Office of Scientific Research under award number FA9550-19-1-0354.

References

- Anselmi, F., Evangelopoulos, G., Rosasco, L., and Poggio, T. Symmetry-adapted representation learning. *Pattern Recognition*, 86:201–208, 2019.
- Benton, G., Finzi, M., Izmailov, P., and Wilson, A. G. Learning invariances in neural networks. *arXiv preprint arXiv:2010.11882*, 2020.
- Bogatskiy, A., Anderson, B., Offermann, J., Roussi, M., Miller, D., and Kondor, R. Lorentz group equivariant neural network for particle physics. In *International Conference on Machine Learning*, pp. 992–1002. PMLR, 2020.
- Bronstein, M. M., Bruna, J., Cohen, T., and Velicković, P. Geometric deep learning: Grids, groups, graphs, geodesics, and gauges. *arXiv preprint arXiv:2104.13478*, 2021.
- Cohen, T. and Welling, M. Learning the irreducible representations of commutative lie groups. In *International Conference on Machine Learning*, pp. 1755–1763, 2014.
- Cohen, T. and Welling, M. Group equivariant convolutional networks. In *International conference on machine learning*, pp. 2990–2999, 2016a.
- Cohen, T., Weiler, M., Kicanaoglu, B., and Welling, M. Gauge equivariant convolutional networks and the icosahedral cnn. In *International Conference on Machine Learning*, pp. 1321–1330. PMLR, 2019a.
- Cohen, T. S. and Welling, M. Steerable cnns. 2016b.

- Cohen, T. S., Geiger, M., Köhler, J., and Welling, M. Spherical cnns. In *International Conference on Learning Representations*, 2018.
- Cohen, T. S., Geiger, M., and Weiler, M. A general theory of equivariant cnns on homogeneous spaces. In *Advances in Neural Information Processing Systems*, pp. 9142–9153, 2019b.
- Cybenko, G. Approximation by superpositions of a sigmoidal function. *Mathematics of control, signals and systems*, 2(4):303–314, 1989.
- Finzi, M., Stanton, S., Izmailov, P., and Wilson, A. G. Generalizing convolutional neural networks for equivariance to lie groups on arbitrary continuous data. *arXiv preprint arXiv:2002.12880*, 2020.
- Finzi, M., Welling, M., and Wilson, A. G. A practical method for constructing equivariant multilayer perceptrons for arbitrary matrix groups. *arXiv preprint arXiv:2104.09459*, 2021.
- Gilmer, J., Schoenholz, S. S., Riley, P. F., Vinyals, O., and Dahl, G. E. Neural message passing for quantum chemistry. In *International Conference on Machine Learning*, pp. 1263–1272. PMLR, 2017.
- Hall, B. *Lie groups, Lie algebras, and representations: an elementary introduction*, volume 222. Springer, 2015.
- Hartford, J., Graham, D., Leyton-Brown, K., and Ravanbakhsh, S. Deep models of interactions across sets. In *International Conference on Machine Learning*, pp. 1909–1918. PMLR, 2018.
- Hornik, K., Stinchcombe, M., and White, H. Multilayer feedforward networks are universal approximators. *Neural networks*, 2(5):359–366, 1989.
- Kipf, T. N. and Welling, M. Semi-supervised classification with graph convolutional networks. *arXiv preprint arXiv:1609.02907*, 2016.
- Kondor, R. and Trivedi, S. On the generalization of equivariance and convolution in neural networks to the action of compact groups. *arXiv preprint arXiv:1802.03690*, 2018.
- Kondor, R., Lin, Z., and Trivedi, S. Clebsch–gordan nets: a fully fourier space spherical convolutional neural network. In *Advances in Neural Information Processing Systems*, pp. 10117–10126, 2018.
- Kunin, D., Sagastuy-Brena, J., Ganguli, S., Yamins, D. L., and Tanaka, H. Neural mechanics: Symmetry and broken conservation laws in deep learning dynamics. *arXiv preprint arXiv:2012.04728*, 2020.
- Landau, L. D. *The classical theory of fields*, volume 2. Elsevier, 2013.
- LeCun, Y., Boser, B., Denker, J. S., Henderson, D., Howard, R. E., Hubbard, W., and Jackel, L. D. Backpropagation applied to handwritten zip code recognition. *Neural computation*, 1(4):541–551, 1989.
- LeCun, Y., Bottou, L., Bengio, Y., and Haffner, P. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998.
- Lee, J. M., Chow, B., Chu, S.-C., Glickenstein, D., Guenther, C., Isenberg, J., Ivey, T., Knopf, D., Lu, P., Luo, F., et al. Manifolds and differential geometry. *Topology*, 643:658, 2009.
- Polyakov, A. M. *Gauge fields and strings*. Routledge, 2018.
- Rao, R. P. and Ruderman, D. L. Learning lie groups for invariant visual perception. In *Advances in neural information processing systems*, pp. 810–816, 1999.
- Ravanbakhsh, S. Universal equivariant multilayer perceptrons. *arXiv preprint arXiv:2002.02912*, 2020.
- Ravanbakhsh, S., Schneider, J., and Póczos, B. Equivariance through parameter-sharing. In *Proceedings of the 34th International Conference on Machine Learning-Volume 70*, pp. 2892–2901. JMLR.org, 2017.

- Schaub, M. T., Benson, A. R., Horn, P., Lippner, G., and Jadbabaie, A. Random walks on simplicial complexes and the normalized hodge 1-laplacian. *SIAM Review*, 62(2):353–391, 2020.
- Sohl-Dickstein, J., Wang, C. M., and Olshausen, B. A. An unsupervised algorithm for learning lie group transformations. *arXiv preprint arXiv:1001.1027*, 2010.
- Wang, R., Walters, R., and Yu, R. Incorporating symmetry into deep dynamics models for improved generalization. In *International Conference on Learning Representations*, 2020.
- Weiler, M. and Cesa, G. General e (2)-equivariant steerable cnns. In *Advances in Neural Information Processing Systems*, pp. 14334–14345, 2019.
- Weiler, M., Geiger, M., Welling, M., Boomsma, W., and Cohen, T. S. 3d steerable cnns: Learning rotationally equivariant features in volumetric data. In *Advances in Neural Information Processing Systems*, pp. 10381–10392, 2018a.
- Weiler, M., Hamprecht, F. A., and Storath, M. Learning steerable filters for rotation equivariant cnns. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 849–858, 2018b.
- Weinberg, S. *The quantum theory of fields*, volume 3. Cambridge university press, 1995.
- Whitaker, E. On the functions which are represented by the expansion of interpolating theory. In *Proc. Roy. Soc. Edinburgh*, volume 35, pp. 181–194, 1915.
- Worrall, D. E., Garbin, S. J., Turmukhambetov, D., and Brostow, G. J. Harmonic networks: Deep translation and rotation equivariance. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 5028–5037, 2017.
- Zaheer, M., Kottur, S., Ravanbakhsh, S., Poczos, B., Salakhutdinov, R. R., and Smola, A. J. Deep sets. In *Advances in neural information processing systems*, pp. 3391–3401, 2017.
- Zhou, A., Knowles, T., and Finn, C. Meta-learning symmetries by reparameterization. *arXiv preprint arXiv:2007.02933*, 2020.

Checklist

The checklist follows the references. Please read the checklist guidelines carefully for information on how to answer these questions. For each question, change the default **[TODO]** to **[Yes]**, **[No]**, or **[N/A]**. You are strongly encouraged to include a **justification to your answer**, either by referencing the appropriate section of your paper or providing a brief inline description. For example:

Did you include the license to the code and datasets? **[Yes]** See Section

Did you include the license to the code and datasets? **[No]** The code and the data are proprietary.

Did you include the license to the code and datasets? **[N/A]**

Please do not modify the questions and only use the provided macros for your answers. Note that the Checklist section does not count towards the page limit. In your paper, please delete this instructions block and only keep the Checklist section heading above along with the questions/answers below.

1. For all authors...
 - (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s contributions and scope? **[Yes]**
 - (b) Did you describe the limitations of your work? **[No]**
 - (c) Did you discuss any potential negative societal impacts of your work? **[N/A]**
 - (d) Have you read the ethics review guidelines and ensured that your paper conforms to them? **[Yes]**
2. If you are including theoretical results...

- (a) Did you state the full set of assumptions of all theoretical results? [Yes]
 - (b) Did you include complete proofs of all theoretical results? [Yes]
3. If you ran experiments...
- (a) Did you include the code, data, and instructions needed to reproduce the main experimental results (either in the supplemental material or as a URL)? [Yes]
 - (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they were chosen)? [Yes]
 - (c) Did you report error bars (e.g., with respect to the random seed after running experiments multiple times)? [Yes]
 - (d) Did you include the total amount of compute and the type of resources used (e.g., type of GPUs, internal cluster, or cloud provider)? [Yes]
4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- (a) If your work uses existing assets, did you cite the creators? [Yes]
 - (b) Did you mention the license of the assets? [N/A]
 - (c) Did you include any new assets either in the supplemental material or as a URL? [Yes]
 - (d) Did you discuss whether and how consent was obtained from people whose data you're using/curating? [N/A]
 - (e) Did you discuss whether the data you are using/curating contains personally identifiable information or offensive content? [N/A]
5. If you used crowdsourcing or conducted research with human subjects...
- (a) Did you include the full text of instructions given to participants and screenshots, if applicable? [N/A]
 - (b) Did you describe any potential participant risks, with links to Institutional Review Board (IRB) approvals, if applicable? [N/A]
 - (c) Did you include the estimated hourly wage paid to participants and the total amount spent on participant compensation? [N/A]

Supplementary Information

A Extended derivations and proofs

Path-ordered Exponential Every element $g \in G$ can be written as a product $g = \prod_a \exp[t_a^i L_i]$ using the matrix exponential (Hall, 2015). This can be done using a path γ connecting I to g on the manifold of G . Here, t_a will be segments of the path γ which add up as vectors to connect I to g . This surjective map can be written as a “path-ordered” (or time-ordered in physics (Weinberg, 1995)) exponential (POE). In the simplest form, POE can be defined by breaking $u = \prod_a \exp[t_a^i L_i]$ down into infinitesimal steps of size $t_a \approx 1/N$ with $N \gg 1$. Choosing γ to be a differentiable path, we can replace the sum over segments $\sum_a t_a$ with an integral along the path $\sum_a t_a = \int_\gamma dt(s) ds$, where $t(s) = d\gamma/ds$ is the tangent vector to the path γ , where $s \in [0, 1]$ parametrizes γ . The POE is then defined as the infinitesimal t_a limit of $g = \prod_a \exp[t_a^i L_i]$. This can be written as

$$\begin{aligned} g &= P \exp \int_\gamma t^i(s) L_i ds = \lim_{N \rightarrow \infty} \prod_{a=1}^N (I + \delta s \gamma^{0i}(s_a) L_i) \\ &= \prod_0^{s_1} ds_0 \gamma^{0i}(s_0) L_i \prod_0^{s_2} ds_1 \gamma^{0j}(s_1) L_j \prod_0^1 ds_N \gamma^{0k}(s_N) L_k \end{aligned} \quad (22)$$

L-conv derivation Let us consider what happens if the kernel in G-conv equation 3 is localized near identity. Let $\kappa_I(u) = c \delta_\eta(u)$, with constants $c \in \mathbb{R}^{m^0 \times m}$ and kernel $\delta_\eta(u) \in \mathbb{R}$ which has support only on an η neighborhood of identity, meaning $\delta_\eta(I + \epsilon^i L_i) \neq 0$ if $|\epsilon^j| < \eta$. This allows us to expand G-conv in the Lie algebra of G to linear order. With $v_\epsilon = I + \epsilon^i L_i$, we have

$$\begin{aligned} [\delta_\eta \star f](g) &= \int_G dv \delta_\eta(v) f(gv) = \int_{k \in k < \eta} dv_\epsilon \delta_\eta(v_\epsilon) f(gv_\epsilon) \\ &= \int d\epsilon \delta_\eta(I + \epsilon^i L_i) f(g(I + \epsilon^i L_i)) \\ &= \int d\epsilon \delta_\eta(I + \epsilon^i L_i) f(g + \epsilon^i g L_i) \\ &= \int d\epsilon \delta_\eta(I + \epsilon^i L_i) \left[f(g) + \epsilon^i g L_i \frac{d}{dg} f(g) + O(\epsilon^2) \right] \\ &= \int d\epsilon \delta_\eta(I + \epsilon^i L_i) \left[I + \epsilon^i g L_i \frac{d}{dg} f(g) + O(\eta^2) \right] \\ &= W^0 \left[I + \bar{\epsilon}^i g L_i \frac{d}{dg} f(g) + O(\eta^2) \right] \end{aligned} \quad (23)$$

where $d\epsilon$ is the integration measure on the Lie algebra induced by the Haar measure dv on G . The $O(\eta^2)$ term arises from integrating the $O(\epsilon^2)$ terms. To see this, note that for an order p function $\phi(\epsilon)$, for $|\epsilon^j| < \eta$, $\int \phi(\epsilon) d\epsilon < \eta^p C$ (for some constant C). Substituting this bound into the integral over the kernel δ_η we get

$$\begin{aligned} &\int d\epsilon \delta_\eta(I + \epsilon^i L_i) \phi(\epsilon) < \int d\epsilon \delta_\eta(I + \epsilon^i L_i) \eta^p C \\ &< \int d\epsilon \delta_\eta(I + \epsilon^i L_i) \eta^p C < \int d\epsilon \delta_\eta(I + \epsilon^i L_i) \eta^p C \end{aligned} \quad (24)$$

In matrix representations, $g L_i \frac{df}{dg} = [g L_i]_\alpha^\beta \frac{df}{dg} = \text{Tr} [g L_i]^T \frac{df}{dg}$. Note that in $g(I + \epsilon^i L_i) \mathbf{x}_0$, the $g L_i \mathbf{x}_0 = \hat{L}_i(g) \mathbf{x}$ come from the pushforward $\hat{L}_i(g) = g L_i g^{-1} \in T_g G$. Here

$$W^0 = c \int d\epsilon \delta_\eta(I + \epsilon^i L_i) \in \mathbb{R}^{m^0 \times m}, \quad \bar{\epsilon}^i = \frac{\int d\epsilon \delta_\eta(I + \epsilon^i L_i) \epsilon^i}{\int d\epsilon \delta_\eta(I + \epsilon^i L_i)} \in \mathbb{R}^m \quad (25)$$

with $k\bar{\epsilon}k < \eta$. When δ_η is normalized, meaning $\int_G \delta_\eta(g) dg = 1$, we have $W^0 = c$ and

$$\bar{\epsilon}^i = \int d\epsilon \delta_\eta(I + \epsilon^i L_i) \epsilon^i$$

Note that with $f(g) \in \mathbb{R}^m$, each $\epsilon^i \in \mathbb{R}^m$ is a matrix. With indices, $f(gv_\epsilon)$ is given by

$$[f(gv_\epsilon)]^a = \int_b f^b(g(\delta_b^a + [\epsilon^i]_b^a L_i)) \quad (26)$$

Similarly, the integration measure $d\epsilon$, which is induced by the Haar measure $dv_\epsilon = d\mu(v_\epsilon)$, is a product $d\epsilon = |J|^{-1} d[\epsilon^i]_b^a$, with $J = \partial v_\epsilon / \partial \epsilon$ being the Jacobian.

Equation 23 is the core of the architecture we are proposing, the Lie algebra convolution or **L-conv**.

L-conv Layer In general, we define Lie algebra convolution (L-conv) as follows

$$\begin{aligned} Q[f](g) &= W^0 \int (I + \bar{\epsilon}^i g L_i) \frac{d}{dg} f(g) \\ &= [W^0]_b f^a \int g \delta_a^b + [\bar{\epsilon}^i]_a^b L_i + O(\bar{\epsilon}^2) \end{aligned} \quad (27)$$

Extended equivariance for L-conv From equation 27 we see that W^0 acts on the output feature indices. Notice that the equivariance of L-conv is due to the way $gv_\epsilon = g(I + \bar{\epsilon}^i L_i)$ appears in the argument, since for $u \in G$

$$u \cdot Q[f](g) = W^0 f(u^{-1} g v_\epsilon) = W^0 [u \cdot f](g v_\epsilon) \quad (28)$$

Because of this, replacing W^0 with a general neural network which acts on the feature indices separately will not affect equivariance. For instance, if we pass L-conv through a neural network to obtain a generalized L-conv Q_σ , we have

$$\begin{aligned} Q_\sigma[f](g) &= \sigma(W f(gv_\epsilon) + b) \\ u \cdot Q_\sigma[f](g) &= Q_\sigma[f](u^{-1} g) = \sigma(W f(u^{-1} g v_\epsilon) + b) \\ &= \sigma(W [u \cdot f](g v_\epsilon) + b) = Q_\sigma[u \cdot f](g) \end{aligned} \quad (29)$$

Thus, L-conv can be followed by any nonlinear neural network as long as it only acts on the feature indices (i.e. a in $f^a(g)$) and not on the spatial indices g in $f(g)$.

A.1 Approximating G-conv using L-conv

Lemma 1 (Approximating the kernel). *Let the kernel $\kappa : G \rightarrow F^0 \rightarrow F$ with $\int_G \kappa(g) k^2 dg < 1$ be continuously differentiable with $k d\kappa(g) / dg k^2 < \xi^2$, and with compact support over $G_0 \subset G$. Let $\kappa_k(g) = c_k \delta_\eta(u_k^{-1} g)$ be a set of kernels with support on an η -neighborhood of $u_k \in G$. Then, $\exists c_k \in F^0 \rightarrow F, u_k \in G$ such that $\tilde{\kappa}$ approximates κ , meaning $\int_G \kappa(g) - \tilde{\kappa}(g) k^2 dg < \zeta^2$ for arbitrary small $\zeta \in \mathbb{R}_+$.*

Proof: The intuition is similar to the universal approximation theorem for neural networks (Hornik et al., 1989; Cybenko, 1989), only generalized to a group manifold instead of \mathbb{R} . Let B_0 be the set of $v_\epsilon = I + \epsilon^i L_i \in \mathfrak{g}$, with $k\epsilon k^2 \leq \eta^2$. Choose a finite set of $u_k \in G$ such that the neighborhoods $B_k = u_k B_0 \subset G$, and such that $\bigcup_k B_k = G_0$. We can show that, for small enough η , the kernel does not change more than ξ over each B_k , allowing us to replace it with a constant localized kernel κ_k with support only on B_k . To see this, consider $g \in B_k$ and $v_\epsilon = I + \epsilon^i L_i \in \mathfrak{g}$, such that $v_\epsilon g \in B_k$. Let γ be a path connecting g to $v_\epsilon g$, via $u(s) = (I + s\epsilon^i L_i)g$, with $s \in [0, 1]$. Using the triangle inequality we have

$$\begin{aligned} k\kappa(g) - \kappa(v_\epsilon g)k^2 &= \int_0^1 ds \left| \frac{d\kappa(u(s))}{ds} \right|^2 \int_0^1 ds \left| \frac{d\kappa(u(s))}{ds} \right|^2 \\ &< \int_0^1 ds \frac{d\kappa}{ds} \xi^2 = k\epsilon k^2 \xi^2 \end{aligned} \quad (30)$$

where $k\epsilon k^2 \ll \eta^2$ because $v_\epsilon g \in B_k$. This means that if we replace $\kappa(g)$ with $\kappa(u_k)$ for any $u_k \in B_k$, our error in approximating κ over B_k is less than $\eta^2 \xi^2$. Setting $\kappa_k(g) = \kappa(u_k) \delta_\eta(u_k^{-1}g)$ with the localized kernels $\delta_\eta(g)$ being any continuously differentiable distribution such that $\int_{B_k} dg \delta_\eta(u_k^{-1}g) = 1$ on all B_k , we get

$$\int_G dg \kappa(g) \int_k \tilde{\kappa}(g) k^2 = \int_k \int_{B_k} dg \kappa(g) \int_k \kappa_k(g) k^2 < \int_k jB_k j \eta^2 \xi^2 = jG_0 j \eta^2 \xi^2 \quad (31)$$

where $jG_0 j$ is the volume of the support of κ

Thus, we can approximate a large class of kernels as $\kappa(g) = \int_k c_k \kappa_k(g)$ where the local kernels $\kappa_k(g) = c_k \delta_\eta(u_k^{-1}g)$ have support only on an η neighborhood of $u_k \in G$. Here $c_k \in \mathbb{R}^m$ are constants and $\delta_\eta(u)$ is as in equation 6. Using this, G-conv equation 3 becomes

$$[\kappa \star f](g) = \int_k c_k \int dv \delta_\eta(u_k^{-1}v) f(gv) = \int_k c_k [\delta_\eta \star f](gu_k). \quad (32)$$

The kernels κ_k are localized around u_k , whereas in L-conv the kernel is around identity. We can compose L-conv layers to move κ_k from u_k to identity.

Lemma 2 (Moving kernels to identity). *The local kernel κ_k can be moved near identity using a multilayer L-conv.*

Proof: In equation 32, write $u_k = v_\epsilon u_k^0$, with $v_\epsilon = I + \epsilon^i L_i \in \mathfrak{g}$. We have

$$\delta_\eta(u_k^{-1}v) = \int_{g \in u_k^0^{-1}v} I + \epsilon^i g L_i \frac{d}{dg} \delta_\eta(g) + O(\epsilon^2) = Q_\epsilon [\delta_\eta](u_k^0^{-1}v) \quad (33)$$

where $Q_\epsilon = I + \epsilon^i \hat{L}_i$ is an L-conv layer with $W^0 = I$. This means that applying one L-conv layer with the parameters above moves the localized kernel along v_ϵ on G . Iterating this further, write u_k as the product of a set of small group elements $u_k = \prod_{a=1}^p v_a$, with $v_a = I + \epsilon^i L_i \in \mathfrak{g}$. Defining L-conv layers $Q_a = I + \epsilon^i \hat{L}_i$, we can write

$$\kappa_k(g) = c_k Q_p \circ Q_1 \delta_\eta(g) \quad (34)$$

meaning κ_k localized around u_k can be written as a p layer L-conv acting on a kernel $\delta_\eta(g)$, localized around the identity of the group. With $k\epsilon_a k < \eta$, the error in u_k is $O(\eta^{p+1})$.

Thus, we conclude that any G-conv equation 3 can be approximated by multilayer L-conv. We can take this result even further, following the main theorem in Kondor & Trivedi (2018), and show that any feedforward equivariant neural network can be approximated using multilayer L-conv with nonlinearities.

Equivariance of nonlinearity Pointwise nonlinearities give equivariant maps between scalar feature maps. To see this, let $\sigma : \mathbb{R} \rightarrow \mathbb{R}$. We extend $\sigma : F \rightarrow F$ by applying σ component-wise. Let $f : S \rightarrow F$ be a scalar feature map (i.e., $g \cdot f(\mathbf{x}) = f(g^{-1}\mathbf{x})$). Then

$$g \cdot (\sigma \circ f)(\mathbf{x}) = \sigma \circ f(g^{-1}\mathbf{x}) = \sigma \circ (g \cdot f)(\mathbf{x}).$$

Since the composition of equivariant maps is equivariant, given equivariant linear mapping $Q : F^S \rightarrow F^S$ (i.e. $g \cdot Q[f] = Q[g \cdot f]$), the layer $f \circ Q$ is equivariant. Hence we have the corollary:

Corollary 2. *Assume G is compact and acts on S transitively. Then any equivariant feedforward neural network (FNN) can be approximated using multilayer L-conv with point-wise nonlinearities. Without the compactness and transitivity hypothesis, multilayer L-conv with pointwise non-linearities can approximate multilayer G-conv with pointwise non-linearities.*

Proof: A FNN is defined as $\sigma_p \circ F_p \circ [\sigma_1 \circ F_1 \circ f](\mathbf{x})$ where F_k are linear and σ_k are point-wise nonlinearities. By Theorem 1 of Kondor & Trivedi (2018), any linear equivariant layer is a G-conv, which by Theorem 1 can be approximated by multilayer L-conv. Therefore, multilayer L-conv with nonlinearity can approximate any equivariant FNN.

Finally, note that to our knowledge it is not known whether every equivariant function can be approximated by equivariant FNN for a Lie group G . Hence, the corollary above is not a universal approximation theorem for equivariant scalar functions in terms of L-conv. However, it does show that multilayer L-conv is equally expressive as other equivariant networks.

Since many datasets such as images deal with discretized spaces, we first need to derive how L-conv acts on such data, discussed next.

A.2 Example of continuous L-conv

The $gL_i \frac{df}{dg}$ in equation 6 can be written in terms of partial derivatives $\partial_\alpha f(\mathbf{x}) = \partial f / \partial \mathbf{x}^\alpha$. In general, using $\mathbf{x}^\rho = g_\rho^\sigma \mathbf{x}_0^\sigma$, we have

$$\frac{df(g\mathbf{x}_0)}{dg_\beta^\alpha} = \frac{d(g_\sigma^\rho \mathbf{x}_0^\sigma)}{dg_\beta^\alpha} \partial_\rho f(\mathbf{x}) = \mathbf{x}_0^\beta \partial_\alpha f(\mathbf{x}) \quad (35)$$

$$gL_i \frac{df}{dg} = [gL_i]_\beta^\alpha \mathbf{x}_0^\beta \partial_\alpha f(\mathbf{x}) = [gL_i \mathbf{x}_0]^\gamma f(\mathbf{x}) = \hat{L}_i f(\mathbf{x}) \quad (36)$$

Hence, for each L_i , the pushforward gL_i generates a flow on S through the vector field $\hat{L}_i = gL_i \frac{d}{dg} = [gL_i \mathbf{x}_0]^\alpha \partial_\alpha$ (Fig. 1). Being a vector field $\hat{L}_i \in TS$ (i.e. 1-tensor), \hat{L}_i is basis independent, meaning for $v \in G$, $\hat{L}_i(v\mathbf{x}) = \hat{L}_i$. Its components transform as $[\hat{L}_i(v\mathbf{x})]^\alpha = [v gL_i \mathbf{x}_0]^\alpha = v_\beta^\alpha \hat{L}_i(\mathbf{x})^\beta$, while the partial transforms as $\partial / \partial [v\mathbf{x}]^\alpha = [v^{-1}]^\gamma_\alpha \partial_\gamma$. Using this relation and Taylor expanding equation 10, we obtain a second form for the group action on L-conv. For $w \in G$, with $\mathbf{y} = w^{-1} \mathbf{x}$ we have

$$Q[f](w^{-1} g\mathbf{x}_0) = W^0 \left(I + \bar{\epsilon}^i [\hat{L}_i]^\alpha [w^{-1}]_\alpha^\beta \frac{\partial}{\partial \mathbf{y}^\beta} \right) f(\mathbf{y})_{\mathbf{y} = w^{-1} \mathbf{x}} \quad (37)$$

1D Translation: Let $G = T_1 = (\mathbb{R}, +)$. A matrix representation for G is found by encoding x as a 2D vector $(x, 1)$. The lift is given by $\mathbf{x}_0 = (0, 1)$ as the origin and $g = \begin{pmatrix} 1 & x \\ 0 & 1 \end{pmatrix}$. The

Lie algebra basis is $L = \begin{pmatrix} 0 & 1 \\ 0 & 0 \end{pmatrix}$. It is easy to check that $gg^0 \mathbf{x}_0 = (x + x^0, 1)$. We also find $gL = L$, meaning L looks the same in all $T_g G$. Close to identity $I = 0$, $v_\epsilon = I + \epsilon L = \epsilon$. We have $gv_\epsilon \mathbf{x}_0 = (g + \epsilon gL) \mathbf{x}_0 = (x + \epsilon, 1)$. Thus, $f(g(I + \epsilon L) \mathbf{x}_0) = f(\mathbf{x}) + \epsilon df(\mathbf{x})/d\mathbf{x}$. This readily generalizes to n D translations T_n (SI A.2.2), yielding $f(\mathbf{x}) + \epsilon^\alpha \partial_\alpha f(\mathbf{x})$.

2D Rotation: Let $G = SO(2)$. The space which $SO(2)$ can lift is not the full \mathbb{R}^2 , but a circle of fixed radius $r = \sqrt{x^2 + y^2}$. Hence we choose $S = S^1$ embedded in \mathbb{R}^2 , with $x = r \cos \theta$ and $y = r \sin \theta$. For the lift, we use the standard 2D representation. We have $\mathbf{x}_0 = (r, 0)$ and (see SI A.2.1)

$$L = \begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix}, \quad g = \exp[\theta L] = \frac{1}{r} \begin{pmatrix} x & y \\ y & x \end{pmatrix}, \quad gL \frac{df}{dg} = (x\partial_y - y\partial_x) f. \quad (38)$$

Physicists will recognize $\hat{L} = (x\partial_y - y\partial_x) = \partial_\theta$ as the angular momentum operator in quantum mechanics and field theories, which generates rotations around the z axis.

Rotation and scaling Let $G = SO(2) \times \mathbb{R}^+$, where the $\mathbb{R}^+ = [0, \infty)$ is scaling. The infinitesimal generator for scaling is identity $L_2 = I$. This group is also Abelian, meaning $[L_2, L] = 0$ ($L \in so(2)$ equation 38). $\mathbb{R}^2/0$ can be lifted to G by choosing $\mathbf{x}_0 = (1, 0)$ in polar coordinates and $\mathbf{x} = g\mathbf{x}_0 = rL_2 \exp[\theta L] \mathbf{x}_0$. We again have $gL \frac{df}{dg} = \partial_\theta f$. We also have $gL_2 = g$, so $gL_2 \mathbf{x}_0 = (x, y)$ and from equation 36, $gL_2 \frac{df}{dg} = (x\partial_x + y\partial_y) f = r\partial_r f$, which is the scaling operation.

A.2.1 Rotation $SO(2)$

With $\mathbf{x}_0 = (r, 0)$

$$g = \begin{pmatrix} \cos \theta & \sin \theta \\ \sin \theta & \cos \theta \end{pmatrix}, \quad = \frac{1}{r} \begin{pmatrix} x & y \\ y & x \end{pmatrix}, \quad L = \begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix}, \quad gL = \frac{1}{r} \begin{pmatrix} y & x \\ x & y \end{pmatrix} \quad (39)$$

$$gL\mathbf{x}_0 = \begin{pmatrix} y \\ x \end{pmatrix} = \begin{pmatrix} \sin \theta \\ \cos \theta \end{pmatrix} \quad (40)$$

To calculate df/dg we note that even after the lift, the function f was defined on S . So we must include the \mathbf{x}_0 in $f(g\mathbf{x}_0)$. Using equation 8, we have

$$\begin{aligned} \frac{df(g\mathbf{x}_0)}{dg} &= \frac{1}{r} \mathbf{x}_0^T \begin{pmatrix} \partial_x f & \partial_y f \\ \partial_y f & \partial_x f \end{pmatrix} = \begin{pmatrix} \partial_x f & \partial_y f \\ 0 & 0 \end{pmatrix} \\ gL \frac{df}{dg} &= \text{Tr} \begin{pmatrix} x\partial_y f & y\partial_x f \\ 0 & x\partial_x f + y\partial_y f \end{pmatrix} = (x\partial_y f \quad y\partial_x f) \end{aligned} \quad (41)$$

A.2.2 Translations T_n

Generalizing the T_1 case, we add a dummy dimension 0 and $\mathbf{x}_0 = (1, 0, \dots, 0)$. The generators are $[L_i]^\nu_\mu = \delta_{i\mu}\delta'_0$ and $g = I + x^i L_i$. Again, $gL_i = L_i + x^j L_j L_i = L_i$ as $L_j L_i = 0$ for all i, j . Hence, $[\hat{L}_i]^\alpha = [gL_i \mathbf{x}_0]^\alpha = \delta_i^\alpha$.

A.3 Group invariant loss

Because G is the symmetry group, f and $g \circ f$ should result in the same optimal parameters. Hence, the minima of the loss function need to be *group invariant*. One way to satisfy this is for the loss itself to be group invariant, which can be constructed by integrating over G (global pooling (Bronstein et al., 2021)). A function $I = \int_G dg F(g)$ is G -invariant as for $w \in G$

$$w \circ I = \int_G w \circ F(g) dg = \int_G F(w^{-1}g) dg = \int_G F(g^\theta) d(wg^\theta) = \int_G F(g^\theta) dg^\theta \quad (42)$$

where we used the invariance of the Haar measure $d(wg^\theta) = dg^\theta$. We can change the integration to $\int_S d^n x$ by change of variable $dg/d\mathbf{x}$. Since we need S to be lifted to G , the lift: $S \rightarrow G$ is injective, the a map $G \rightarrow S$ need not be. S is homeomorphic to G/H , where $H \subset G$ is the stabilizer of the origin, i.e. $h\mathbf{x}_0 = \mathbf{x}_0, \forall h \in H$. Since $F(g\mathbf{x}_0) = F(hg\mathbf{x}_0)$, we have

$$I = \int_G F(g) dg = \int_H \int_{G/H} dh \int_{G/H} dg^\theta F(g^\theta) = V_H \int_{G/H} dg^\theta F(g^\theta) \quad (43)$$

Since $G/H \cong S$, the volume forms $dg^\theta = V_H d^n x$ can be matched for some parametrization.

A.3.1 MSE Loss

The MSE is given by $I = \int_G dg kQ[f_n](g)k^2$, where f_n are data samples and $Q[f]$ is L-conv or another G -equivariant function. In supervised learning the input is a pair f_n, y_n . G can also act on the labels y_n . We assume that y_n are either also scalar features $y_n : S \rightarrow \mathbb{R}^{m_y}$ with a group action $g \circ y_n(\mathbf{x}) = y_n(g^{-1}\mathbf{x})$ (e.g. f_n and y_n are both images), or that y_n are categorical. In the latter case $g \circ y_n = y_n$ because the only representations of a continuous G on a discrete set are constant. We can concatenate the inputs to $\phi_n = [f_n | y_n]$ with a well-defined G action $g \circ \phi_n = [g \circ f_n | g \circ y_n]$. The collection of combined inputs $\Phi = (\phi_1, \dots, \phi_N)^T$ is an $(m + m_y) \times N$ matrix. Using equations 6 and 8, the MSE loss with parameters $W = \bar{f}W^0, \bar{e}g$ becomes

$$\begin{aligned} I[\Phi; W] &= \int_G dg L[\Phi; W] = \int_G dg W^0 \left(I + \bar{\epsilon}^i [\hat{L}_i]^\alpha \partial_\alpha \right) \Phi(g)^2 \\ &= 2 \int_G dg kW^0 \Phi k^2 + W^i [\hat{L}_i]^\alpha \partial_\alpha \Phi^2 + 2\Phi^T W^{0T} W^i [\hat{L}_i]^\alpha \partial_\alpha \Phi \end{aligned} \quad (44)$$

$$= \int_S \frac{d^n x}{\frac{\partial x}{\partial g}} \Phi^T \mathbf{m}_2 \Phi + \partial_\alpha \Phi^T \mathbf{h}^{\alpha\beta} \partial_\beta \Phi + [\hat{L}_i]^\alpha \partial_\alpha \Phi^T \mathbf{v}^i \Phi \quad (45)$$

where $\frac{\partial x}{\partial g}$ is the determinant of the Jacobian, $W^i = W^0 \bar{\epsilon}^i$ and

$$\mathbf{m}_2 = W^{0T} W^0, \quad \mathbf{h}^{\alpha\beta}(\mathbf{x}) = \bar{\epsilon}^{iT} \mathbf{m}_2 \bar{\epsilon}^j [\hat{L}_i]^\alpha [\hat{L}_j]^\beta, \quad \mathbf{v}^i = \mathbf{m}_2 \bar{\epsilon}^i. \quad (46)$$

From equation 44 to 45 we used the fact that W^0 and W^i do not depend on \mathbf{x} (or g) to write

$$2\Phi^T W^{0T} W^i [\hat{L}_i]^\alpha \partial_\alpha \Phi = [\hat{L}_i]^\alpha \partial_\alpha \Phi^T W^{0T} W^i \Phi = [\hat{L}_i]^\alpha \partial_\alpha \Phi^T \mathbf{m}_2 \bar{\epsilon}^i \Phi \quad (47)$$

Note that \mathbf{h} has feature space indices via $[\bar{\epsilon}^{iT} \mathbf{m}_2 \bar{\epsilon}^j]_{ab}$, with index symmetry $\mathbf{h}_{ab}^{\alpha\beta} = \mathbf{h}_{ba}^{\beta\alpha}$. When $F = \mathbb{R}$ (i.e. f is a 1D scalar), $\mathbf{h}^{\alpha\beta}$ becomes a Riemannian metric for S . In general \mathbf{h} combines a 2-tensor $\mathbf{h}_{ab} = \mathbf{h}_{ab}^{\alpha\beta} \partial_\alpha \partial_\beta$ $\mathbb{Z}TS \quad TS$ with an inner product $h^T \mathbf{h}^{\alpha\beta} f$ on the feature space F . Hence $\mathbf{h} \mathbb{Z}TS \quad TS \quad F \quad F$ is a $(2, 2)$ -tensor, with F being the dual space of F .

Loss invariant metric transformation The metric \mathbf{h} transforms equivariantly as a 2-tensor. As discussed under equation 8, $[\hat{L}_i(v\mathbf{x})]^\alpha = v_\beta^\alpha \hat{L}_i(\mathbf{x})^\beta$ and

$$v \mathbf{h}^{\alpha\beta} = \mathbf{h}^{\alpha\beta}(v^{-1}\mathbf{x}) = [v^{-1}]_\rho^\alpha [v^{-1}]_\gamma^\beta \mathbf{h}^{\rho\gamma}(\mathbf{x}), \quad (v \in G). \quad (48)$$

Note that $v \mathbf{m}_2 = \mathbf{m}_2$ since f_n and y_n are scalars. For example, let $G = SO(2)$ and $R(\xi) \in SO(2)$ be rotation by angle ξ . Since there is only one $L_i = L$, the metric factorizes to

$$\mathbf{h}_{ab}^{\alpha\beta} = [\bar{\epsilon}^T \mathbf{m}_2 \bar{\epsilon}]_{ab} \quad [\hat{L} \hat{L}^T]^{\alpha\beta} \quad (49)$$

To find $R(\xi) \mathbf{h}$ we only need to calculate $R(\xi)^{-1} \hat{L}$. With $g = R(\theta)$, we have $\hat{L}(\mathbf{x}) = R(\theta) L \mathbf{x}_0 = (y, x) = r(\sin \theta, \cos \theta)$ from equation 40. Therefore, $R(\xi)^{-1} \hat{L} = R(\theta - \xi) \hat{L} = \hat{L}(R(\xi^{-1})\mathbf{x})$. Using equation 19 in equation 48, the transformed metric becomes

$$R(\xi) \mathbf{h}^{\alpha\beta}(R(\theta)\mathbf{x}_0) = \bar{\epsilon}^T \mathbf{m}_2 \bar{\epsilon} \quad [R(\theta - \xi) \hat{L}]^\alpha [R(\theta - \xi) \hat{L}]^\beta = \mathbf{h}^{\alpha\beta}(R(\theta - \xi)\mathbf{x}_0), \quad (50)$$

A.3.2 Third term as a boundary term

Since terms in equation 18 are scalars, they can be evaluated in any basis. If S can be lifted to multiple Lie groups, either group can be used to evaluate equation 18. For example $\mathbb{R}^n/0$ can be lifted to both T_n and $SO(n) \times \mathbb{R}_+$. For the translation group $G = T_n$ we have $gL_i = L_i$ and $[\hat{L}_i]^\alpha = \delta_i^\alpha$ (SI A.2.2) and $\Phi^j \partial g / \partial x^j = 1$ and $dg = d^n x$. Thus, the last term in equation 18 simplifies to a complete divergence $d^n x \partial_i (\Phi^T \mathbf{v}^i \Phi)$. Using the generalized Stoke's theorem $\int_S dw = \int_{\partial S} w$, the last term in equation 18 becomes a boundary term. When S is non-compact, the last term is $I_\partial = \int_{\partial S} d\Sigma_i \Phi^T \mathbf{v}^i \Phi$, where $d\Sigma_i$ is the normal times the volume form of the $(n - 1)$ D boundary ∂S and is in the radial direction (e.g. for $S = \mathbb{R}^n$ the boundary is a hyper-sphere $\partial S = S^{n-1}$). Generally we expect the features ϕ_n to be concentrated in a finite region of the space and that they go to zero as $r \rightarrow \infty$ (if they don't the loss term $\Phi^T \mathbf{m}_2 \Phi$ will diverge). Thus, the last term in equation 18 generally becomes a vanishing boundary term and does not matter.

A.3.3 MSE Loss for translation group T_n

We have $gL_i = L_i$ and $[\hat{L}_i]^\alpha = \delta_i^\alpha$ (SI A.2.2), the last term in equation 18 becomes a complete divergence $I_\partial = \int_{\partial S} d^n x \partial_i (\Phi^T \mathbf{v}^i \Phi)$. Using the generalized Stoke's theorem $\int_S dw = \int_{\partial S} w$, when S is non-compact, $I_\partial = \int_{\partial S} d\Sigma_i \Phi^T \mathbf{v}^i \Phi$. Here $d\Sigma_i$ is the normal times the volume form of the $(n - 1)$ D boundary ∂S (e.g. for $S = \mathbb{R}^n$ the boundary is a hyper-sphere $\partial S = S^{n-1}$). Generally we expect the features ϕ_n to be concentrated in a finite region of the space and that they go to zero as $r \rightarrow \infty$ (if they don't the loss term $\Phi^T \mathbf{m}_2 \Phi$ will diverge). Thus, the last term in equation 18 generally becomes a vanishing boundary term and does not matter.

Next, the second term in equation 18 can be worked out as

$$\hat{L}_i^\alpha \partial_\alpha \phi^T \bar{\epsilon}^{iT} \mathbf{m}_2 \bar{\epsilon}^j \hat{L}_j^\beta \partial_\beta \phi = \partial_j \phi^T \bar{\epsilon}^{iT} \mathbf{m}_2 \bar{\epsilon}^j \partial_i \phi = \partial_j \phi^T \mathbf{h}^{ji} \partial_i \phi \quad (51)$$

where $\mathbf{h}^{ji} = \bar{\epsilon}^{iT} \mathbf{m}_2 \bar{\epsilon}^j$ is a general, space-independent metric compatible with translation symmetry. When the weights $[W^i]_b^a \sim \mathcal{N}(0, 1)$ are random Gaussian, we have $W^{jT} W^i \sim m^2 \delta^{ij}$ and we recover the Euclidean metric. With the 1st term vanishing, the loss function equation 18 has a striking resemblance to a Lagrangian used in physics, as we discuss next.

A.3.4 Boundary term with spherical symmetry

When $S = \mathbb{R}^n$ and $G = T_n$, the third term becomes a boundary term. But we can also have $G = SO(n) \times \mathbb{R}_+$ (spherical symmetry and scaling). The boundary $\partial S = S^{n-1}$, which has an $SO(n)$ symmetry. The normal $d\Sigma(\mathbf{x})$ is a vector pointing in the radial direction and g is the lift for \mathbf{x} . Since $g \in SO(n)$, we have

$$d\Sigma_\beta [gL_\alpha \mathbf{x}_0]^\beta = d\Sigma^T g L_\alpha \mathbf{x}_0 = [g^T d\Sigma]^T L_\alpha \mathbf{x}_0 \quad (52)$$

Since $g \in SO(n)$, $g^T = g^{-1}$ and $g^T d\Sigma(g\mathbf{x}_0) = d\Sigma(\mathbf{x}_0) = V_{n-1}\mathbf{x}_0$, meaning the normal vector is rotated back toward \mathbf{x}_0 . Here V_{n-1} is the volume of the boundary S^{n-1} . Hence we have

$$d\Sigma^T g L_i \mathbf{x}_0 = \mathbf{x}_0^T L_i \mathbf{x}_0 = 0 \quad (53)$$

for all generators $L_i \in so(n)$ because $L_i = -L_i^T$ and hence diagonal entries like $\mathbf{x}_0^T L_i \mathbf{x}_0$ are zero. Only the scaling generator $L_0 = I$ we have $\mathbf{x}_0^T L_0 \mathbf{x}_0 = 1$. This means that the last term in equation 18 can be nonzero at the boundary only if $\Phi \mathbf{v}^i \Phi$ is in the radial direction, meaning $\bar{\epsilon}^0 = 0$, and Φ does not vanish at the boundary. However, a non-vanishing Φ at the boundary results in diverging loss unless the mass matrix \mathbf{m}_2 has eigenvalues equal to zero. This is what happens relativistic theories where light rays can have nonzero Φ at infinity because they are massless.

A.4 Robustness to random noise

Equivariant neural networks are hoped to be more robust than others. One way to check this is to see how the network would perform for an input $\phi^\theta = \phi + \delta\phi$ which adds a small perturbation $\delta\phi$ to a real data point ϕ . Robustness to such perturbation would mean that, for optimal parameters W , the loss function would not change, i.e. $I[\phi^\theta; W] = I[\phi; W]$. This can be cast as a variational equation, requiring I to be minimized around real data points ϕ . Writing $I[\phi; W] = \int_S d^n x L[\phi; W]$, we have

$$\delta I[\phi; W] = \int_S d^n x \frac{\partial L}{\partial \phi^a} \delta \phi^a + \frac{\partial L}{\partial (\partial_\alpha \phi^a)} \partial_\alpha (\delta \phi^a) \quad (54)$$

Doing a partial integration on the second term, we get

$$\begin{aligned} \delta I[\phi; W] &= \int_S d^n x \frac{\partial L}{\partial \phi^b} \delta \phi^b + \int_S d^n x \partial_\alpha \frac{\partial L}{\partial (\partial_\alpha \phi^b)} \delta \phi^b \\ &= \int_S d^n x \frac{\partial L}{\partial \phi^b} \delta \phi^b + \int_{\partial S} d^{n-1} \Sigma_\alpha \frac{\partial L}{\partial (\partial_\alpha \phi^b)} \delta \phi^b \end{aligned} \quad (55)$$

If we want equivariant networks to be robust, then both terms in equation 55 need to be zero. We show that the first term is the classic Euler-Lagrange (EL) equation, and the second term is related to conservation laws.

We use the Stoke's theorem to change the second term to a boundary integral. Since features ϕ have finite support, $\phi(\mathbf{x}) \rightarrow 0$ as $|\mathbf{x}| \rightarrow \infty$, and the boundary term vanishes. The first term in equation 55 is the classic Euler-Lagrange (EL) equation. Thus, requiring robustness, i.e. $\delta I[\phi; W]/\delta \phi = 0$ means for optimal parameters W , the data ϕ satisfies the EL equations

$$\text{Robustness to random noise () EL: } \frac{\partial L}{\partial \phi^b} - \partial_\alpha \frac{\partial L}{\partial (\partial_\alpha \phi^b)} = 0 \quad (56)$$

Applying this to the MSE loss equation 18, equation 56 becomes

$$\mathbf{m}_2 \phi - \partial_\alpha (jJ \mathbf{h}^{\alpha\beta} \partial_\beta \phi) - \partial_\alpha (jJ \mathbf{v}^i [\hat{L}_i]^\alpha \phi) = 0 \quad (57)$$

where $jJ = |j\partial g/\partial x^j|$ is the determinant of the Jacobian. For the translation group, equation 57 becomes a Helmholtz equation

$$\mathbf{h}^{ij} \partial_i \partial_j \phi = \bar{\epsilon}^i \mathbf{m}_2 \bar{\epsilon}^j \partial_i \partial_j \phi = \mathbf{m}_2 \phi \quad (58)$$

where $\mathbf{h}^{ij} \partial_i \partial_j = \nabla^2$ is the Laplace-Beltrami operator with \mathbf{h} as the metric.

A.5 Conservation laws

The equivariance condition equation 2 can be written for the integrand of the loss $L[\phi; W]$. Since G is the symmetry of the system, transforming an input $\phi \rightarrow w \phi$ by $w \in G$ the integrand changes equivariantly as $L[w \phi] = w L[\phi]$. Now, let w be an infinitesimal $w = I + \eta^i L_i$. The action $w \phi$ can be written as a Taylor expansion, similar to the one in L-conv, yielding

$$\begin{aligned} w \phi(\mathbf{x}) &= \phi(w^{-1} \mathbf{x}) = \phi((I - \eta^i L_i) \mathbf{x}) = \phi(\mathbf{x}) - \eta^i [L_i \mathbf{x}]^\alpha \partial_\alpha \phi(\mathbf{x}) \\ &= \phi(\mathbf{x}) + \delta \mathbf{x}^\alpha \partial_\alpha \phi(\mathbf{x}) = \phi(\mathbf{x}) + \delta \phi(\mathbf{x}) \end{aligned} \quad (59)$$

with $\delta \mathbf{x}^\alpha = \eta^i [L_i \mathbf{x}]^\alpha$ and $\delta \phi = \delta \mathbf{x}^\alpha \partial_\alpha \phi$. Similarly, we have $w L = L + \delta \mathbf{x}^\alpha \partial_\alpha L$. Next, we can use the chain rule to calculate $L[w \phi]$.

$$\begin{aligned} L[w \phi] &= L[\phi(\mathbf{x}) + \delta \phi(\mathbf{x})] = L[\phi] + \frac{\partial L}{\partial \phi^b} \delta \phi^b + \frac{\partial L}{\partial (\partial_\alpha \phi^b)} \delta \partial_\alpha \phi^b \\ &= L[\phi] + \frac{\partial L}{\partial \phi^b} \left(\partial_\alpha \frac{\partial L}{\partial (\partial_\alpha \phi^b)} \delta \phi^b + \partial_\alpha \frac{\partial L}{\partial (\partial_\alpha \phi^b)} \delta \phi^b \right) \end{aligned} \quad (60)$$

where we used the fact that $\delta \mathbf{x} = \eta^i L_i \mathbf{x}$ can vary independently from \mathbf{x} (because of η^i), and so $\delta \partial_\alpha \phi^b = \partial_\alpha \delta \phi^b$. The same way, $\delta \mathbf{x}^\alpha \partial_\alpha L = \partial_\alpha (\delta \mathbf{x}^\alpha L)$. Now, if ϕ are the real data and the parameters in L minimize generalization error, then L satisfies equation 57. This means that the first term in equation 60 vanishes. Setting the second term equal to $w L$ we get

$$L[w \phi] - w L[\phi] = \partial_\alpha \frac{\partial L}{\partial (\partial_\alpha \phi^b)} \delta \phi^b - \delta \mathbf{x}^\alpha L = 0 \quad (61)$$

Thus, the terms in the brackets are divergence free. These terms are called a Noether conserved current J^α . In summary

$$\text{Noether current: } J^\alpha = \frac{\partial L}{\partial (\partial_\alpha \phi^b)} \delta \phi^b - \frac{\partial L}{\partial \mathbf{x}^\alpha} \delta \mathbf{x}^\alpha, \quad \delta I[\phi; W] = 0 \implies \partial_\alpha J^\alpha = 0 \quad (62)$$

J captures the change of the Lagrangian L along symmetry direction \hat{L}_i . Plugging $\delta \phi = \delta \mathbf{x}^\alpha \partial_\alpha \phi$ from equation 59 we find

$$\partial_\alpha \frac{\partial L}{\partial (\partial_\alpha \phi^b)} \delta \phi^b - \delta \mathbf{x}^\alpha L = \delta \mathbf{x}^\beta \partial_\alpha \frac{\partial L}{\partial (\partial_\alpha \phi^b)} \partial_\beta \phi^b - \delta_\beta^\alpha L = \delta \mathbf{x}^\beta \partial_\alpha T_\beta^\alpha. \quad (63)$$

T_β^α is known as the stress-energy tensor in physics (Landau, 2013). It is the Noether current associated with space (or space-time) variations $\delta \mathbf{x}$. It appears here because G acts on the space, as opposed to acting on feature dimensions. For the MSE loss we have

$$T_\beta^\alpha = \frac{\partial L}{\partial (\partial_\alpha \phi^b)} \partial_\beta \phi^b - \delta_\beta^\alpha L = \partial_\rho \phi^T \delta_\beta^\lambda \mathbf{h}^{\alpha\rho} - \delta_\beta^\alpha \mathbf{h}^{\rho\lambda} \partial_\lambda \phi = \phi^T \mathbf{m}_2 \phi \quad (64)$$

It would be interesting to see if the conserved currents can be used in practice as an alternative way for identifying or discovering symmetries.

B Tensor notation details

If the dataset being analyzed is in the form of $f(\mathbf{x})$ for some sample of points \mathbf{x} , together with derivatives $\nabla f(\mathbf{x})$, we can use the L-conv formulation above. However, in many datasets, such as images, $f(\mathbf{x})$ is given as a finite dimensional array or tensor, with \mathbf{x} taking values over a grid. Even though the space S is now discrete, the group which acts on it can still be continuous (e.g. image rotations). Let $S = \{f_{\mathbf{x}_0}, \dots, f_{\mathbf{x}_{d-1}}\}$ contain d points. Each \mathbf{x}_μ represents a coordinate in higher dimensional grid. For instance, on a 10×10 image, \mathbf{x}_0 is $(x, y) = (0, 0)$ point and \mathbf{x}_{99} is $(x, y) = (9, 9)$.

Feature maps To define features $f(\mathbf{x}_\mu) \in \mathbb{R}^m$ for $\mathbf{x}_\mu \in S$, we embed $\mathbf{x}_\mu \in \mathbb{R}^d$ and encode them as the canonical basis (one-hot) vectors with components $[\mathbf{x}_\mu]^\nu = \delta_\mu^\nu$ (Kronecker delta), e.g. $\mathbf{x}_0 = (1, 0, \dots, 0)$. The feature space becomes $F = \mathbb{R}^d \times \mathbb{R}^m$, meaning feature maps $\mathbf{f} \in F$ are $d \times m$ tensors, with $f(\mathbf{x}_\mu) = \mathbf{x}_\mu^T \mathbf{f} = \mathbf{f}_\mu$.

Group action Any subgroup $G \subseteq \text{GL}_d(\mathbb{R})$ of the general linear group (invertible $d \times d$ matrices) acts on \mathbb{R}^d and F . Since $\mathbf{x}_\mu \in \mathbb{R}^d$, $g \in G$ also naturally act on \mathbf{x}_μ . The resulting $\mathbf{y} = g \mathbf{x}_\mu$ is a linear combination $\mathbf{y} = c^\nu \mathbf{x}_\nu$ of elements of the discrete S , not a single element. The action of G on \mathbf{f} and \mathbf{x} , can be defined in multiple equivalent ways. We define $f(g \mathbf{x}_\mu) = \mathbf{x}_\mu^T g^T \mathbf{f}$, $g \in G$. For $w \in G$ we have

$$w f(\mathbf{x}_\mu) = f(w^{-1} \mathbf{x}_\mu) = \mathbf{x}_\mu^T w^{-1T} \mathbf{f} = [w^{-1} \mathbf{x}]^T \mathbf{f} \quad (65)$$

Dropping the position \mathbf{x}_μ , the transformed features are matrix product $w \mathbf{f} = w^{-1T} \mathbf{f}$.

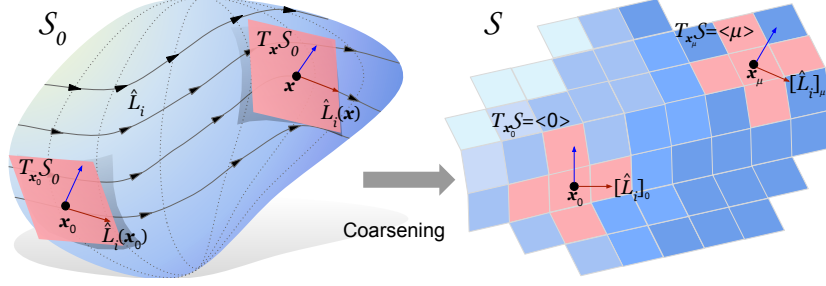


Figure 5: **Manifold vs. discretized Space** While real systems can have a continuous manifold S_0 as their base space, often the data collected from is a discrete array S . The discretization (coarsening) will induce some of the topology of S_0 on S as a graph. Graph neighborhoods $\langle \mu \rangle$ on the discrete S represent tangent spaces $T_x S$ and approximate $T_x S_0$. The lift takes $\mathbf{x} \in S_0$ to $g \in G$, and maps the tangent spaces $T_x S_0 \rightarrow T_x G$. Each Lie algebra basis $L_i \in \mathfrak{g} = T_x G$ generates a vector field on the tangent bundle TG via the pushforward as $L^{(g)} = gL_i g^{-1}$. Due to the lift, each L_i also generates a vector field $\hat{L}_i^\alpha(\mathbf{x})\partial_\alpha = [gL_i \mathbf{x}_0]^\alpha \partial_\alpha$, with $\mathbf{x} = g\mathbf{x}_0$. Analogously, on S we get a vector field $[\hat{L}_i]_\mu = g_\mu L_i \mathbf{x}_0$ on TS , with $\mathbf{x}_\mu = g_\mu \mathbf{x}_0$. Note that depicting S and S_0 as 2D is only for convenience. They may have any dimensions.

G-conv and L-conv in tensor notation Writing G-conv equation 3 in the tensor notation we have

$$[\kappa \star f](g\mathbf{x}_0) = \int_G \kappa(v) f(gv\mathbf{x}_0) dv = \mathbf{x}_0^T \int_G v^T g^T \mathbf{f} \kappa^T(v) dv = \mathbf{x}_0^T [\mathbf{f} \star \kappa](g) \quad (66)$$

where we moved $\kappa^T(v) \in \mathbb{R}^m \rightarrow \mathbb{R}^{m^0}$ to the right of \mathbf{f} because it acts as a matrix on the output index of \mathbf{f} . The equivariance of equation 66 is readily checked with $w \in G$

$$w [\mathbf{f} \star \kappa](g) = [\mathbf{f} \star \kappa](w^{-1}g) = \int_G v^T g^T w^{-1T} \mathbf{f} \kappa^T(v) dv = [(w^{-1} \mathbf{f}) \star \kappa](g) \quad (67)$$

where we used $[w^{-1}g]^T \mathbf{f} = g^T w^{-1T} \mathbf{f}$. Similarly, we can rewrite L-conv equation 6 in the tensor notation. Defining $v_\epsilon = I + \bar{\epsilon}^i L_i$

$$\begin{aligned} Q[\mathbf{f}](g) &= W^0 f \cdot g \cdot I + \bar{\epsilon}^i L_i = \mathbf{x}_0^T [I + \bar{\epsilon}^i L_i]^T g^T \mathbf{f} W^{0T} \\ &= \mathbf{x}_0 + \bar{\epsilon}^i [gL_i \mathbf{x}_0]^T \mathbf{f} W^{0T}. \end{aligned} \quad (68)$$

Here, $\hat{L}_i = gL_i \mathbf{x}_0$ is exactly the tensor analogue of pushforward vector field \hat{L}_i in equation 8. We will make this analogy more precise below. The equivariance of L-conv in tensor notation is again evident from the $g^T \mathbf{f}$, resulting in

$$Q[w^{-1} \mathbf{f}](g) = \mathbf{x}_0^T v_\epsilon^T g^T w^{-1T} \mathbf{f} W^{0T} = Q[\mathbf{f}](w^{-1}g) = w^{-1} Q[\mathbf{f}](g) \quad (69)$$

Next, we will discuss how to implement equation 68 in practice and how to learn symmetries with L-conv. We will also discuss the relation between L-conv and other neural architectures.

B.1 Constraints from topology on tensor L-conv

To implement equation 68 we need to specify the lift and the form of L_i . We will now discuss the mathematical details leading to an easy to implement form of equation 68.

Topology Although the discrete space S is a set of points, in many cases it has a topology. For instance, S can be a discretization of a manifold S_0 , or vertices on a lattice or a general graph. We encode this topology in an undirected graph (i.e. 1-simplex) with vertex set S edge set E . Instead of the commonly used graph adjacency matrix \mathbf{A} , we will use the incidence matrix $\mathbf{B} : S \rightarrow E \rightarrow \mathbb{R}, 1, -1, g$. $\mathbf{B}_\alpha^\mu = 1$ or -1 if edge α starts or ends at node μ , respectively, and $\mathbf{B}_\alpha^\mu = 0$ otherwise (undirected graphs have pairs of incoming and outgoing edges). Similar to the continuous case we will denote the topological space (S, E, \mathbf{B}) simply by S .

Figure 5 summarizes some of the aspects of the discretization as well as analogies between S_0 and S . Technically, the group G_0 acting on S_0 and G acting on S are different. But we can find a group G which closely approximates G_0 (see SI B.2). For instance, Rao & Ruderman (1999) used Shannon-Whittaker interpolation theorem (Whittaker, 1915) to define continuous 1D translation and 2D rotation groups on discrete data. We return to this when expressing CNN as L-conv in 5.

Neighborhoods as discrete tangent bundle \mathcal{B} is useful for extending differential geometry to graphs (Schaub et al., 2020). Define the neighborhood $\langle \mu \rangle = \alpha \in E \mid \mathcal{B}_\alpha^\mu = 1$ of \mathbf{x}_μ as the set of outgoing edges. $\langle \mu \rangle$ can be identified with $T_{\mathbf{x}_\mu} S$ as for $\alpha \in \langle \mu \rangle$, $\mathcal{B}_\alpha \mathbf{f} = \mathbf{f}_\mu - \mathbf{f}_\nu - \partial_\alpha \mathbf{f}$, where μ and ν are the endpoints of edge α . This relation becomes exact when S is an n D square lattice with infinitesimal lattice spacing. The set of all neighborhoods is \mathcal{B} itself and encodes the approximate tangent bundle TS . For some operator $C : S \rightarrow F \rightarrow F$ acting on \mathbf{f} we will say $C \langle \mu \rangle$ if its action remains within vertices connected to μ , meaning

$$C \langle \mu \rangle : \quad C \mathbf{f} = \sum_{\alpha \in \langle \mu \rangle} \tilde{C}^\alpha \mathcal{B}_\alpha^\nu \mathbf{f}_\nu \quad (70)$$

Lift and group action Lie algebra elements L_i by definition take the origin \mathbf{x}_0 to points close to it. Thus, for small enough η , $(I + \eta L_i) \mathbf{x}_0 \in \langle 0 \rangle$ and so $[L_i \mathbf{x}_0]^T \mathbf{f} = [\hat{L}_i]_0^\rho \mathbf{f}_\rho = \sum_{\alpha \in \langle 0 \rangle} [\hat{\ell}_i]_0^\alpha \mathcal{B}_\alpha^\rho \mathbf{f}_\rho$. The coefficients $[\hat{\ell}_i]_0^\alpha \in \mathbb{R}$ are in fact the discrete version of \hat{L}_i components from equation 8. For the pushforward $g L_i g^{-1}$, we define the lift via $\mathbf{x}_\mu = g_\mu \mathbf{x}_0$. We require the G -action to preserve the topology of S , meaning points which are close remain close after the G -action. As a result, $\langle \mu \rangle$ can be reached by pushing forward elements in $\langle 0 \rangle$. Thus, for each i , $\exists \eta^{-1}$ such that $g_\mu (I + \eta L_i) \mathbf{x}_0 \in \langle \mu \rangle$, meaning for a set of coefficients $[\hat{L}_i]_\mu^\nu \in \mathbb{R}$ we have

$$[g_\mu (I + \eta L_i) \mathbf{x}_0]^T \mathbf{f} = \mathbf{f}_\mu + \eta \sum_{\alpha \in \langle \mu \rangle} [\hat{\ell}_i]_\mu^\alpha \mathcal{B}_\alpha^\nu \mathbf{f}_\nu \quad (71)$$

where $\mathbf{f}_\mu = \mathbf{x}_\mu^T \mathbf{f}$. Acting with $[g_\mu L_i \mathbf{x}_0]^T \mathbf{x}_\nu$ and inserting $I = \sum_\rho \mathbf{x}_\rho \mathbf{x}_\rho^T$ we have

$$\begin{aligned} [\hat{L}_i]_\mu^\nu &= [\hat{\ell}_i]_\mu^\alpha \mathcal{B}_\alpha^\nu = [g_\mu L_i \mathbf{x}_0]^\nu = \sum_\rho [g_\mu \mathbf{x}_\rho \mathbf{x}_\rho^T L_i \mathbf{x}_0]^T \mathbf{x}_\nu \\ &= \sum_{\rho \in \langle 0 \rangle} [L_i]_0^\rho [\mathbf{x}_\nu^T g_\mu \mathbf{x}_\rho]^T = [L_i]_0^\rho [g_\mu]_\rho^\nu = [\hat{\ell}_i]_0^\alpha \mathcal{B}_\alpha^\rho [g_\mu]_\rho^\nu \end{aligned} \quad (72)$$

This $\hat{L}_i = \sum_\alpha \hat{\ell}_i^\alpha \mathcal{B}_\alpha$ is the discrete S version of the vector field $\hat{L}_i(\mathbf{x}) = [g L_i \mathbf{x}_0]^\alpha \partial_\alpha$ in equation 8.

B.2 Approximating a symmetry and discretization error

Discretization error While systems such as crystalline solids are discrete in nature, many other datasets such as images result from discretization of continuous data. The discretization (or ‘‘coarsening’’ (Bronstein et al., 2021)) will modify the groups that can act on the space. For example, first rotating a shape by $SO(2)$ then taking a picture is different from first taking a picture then rotating the picture (i.e. group action and discretization do not commute). Nevertheless, in most cases in physics and machine learning the symmetry group G_0 of the space before discretization has a small Lie algebra dimension n , (e.g. $SO(3)$, $SE(3)$, $SO(3, 1)$ etc). Usually the resolution of the discretization is $d \ll n$. In this case, there always exist some $G \subset GL_d(\mathbb{R})$ which approximates G_0 reasonably well. The approximation means $\exists g_0 \in G_0, \exists g \in G$ such that the error $L_G = kg_0 - f(\mathbf{x}_\mu) - \mathbf{x}_\mu^T g \mathbf{f} k^2 < \eta^2$ where η depends on the resolution of the discretization. Minimizing the error L_G can be the process of identifying the G which best approximates G_0 . We will denote this approximate similarity as $G \approx G_0$. For example, Rao & Ruderman (1999) used the Shannon-Whittaker Interpolation theorem (Whittaker, 1915) to translate discrete 1D signals (features) by arbitrary, continuous amounts. In this case the transformed features are $\mathbf{f}_\mu^0 = g(z)_\mu^\nu \mathbf{f}_\nu$, where $g(z)_\mu^\nu = \frac{1}{d} \prod_{p=1}^{d/2} \cos \frac{2\pi p}{d} (z + \mu - \nu)$ approximates the shift operator for continuous z . The $g(z)$ form a group because $g(w)g(z) = g(w + z)$, which is a representation for periodic 1D shifts. Rao & Ruderman (1999) also use a 2D version of the interpolation theorem to approximate $SO(2)$. In practice, we can assume the true symmetry to be G , as we only have access to the discretized data and can’t measure G_0 directly.

B.3 Comparison with other symmetry discovery methods

Meta-learning Symmetries by Reparameterization Recently Zhou et al. (2020) also introduced an architecture which can learn equivariances from data. We would like to highlight the differences between their approach and ours, specifically Proposition 1 in Zhou et al. (2020). Assuming a discrete group $G = \{g_1, \dots, g_n\}$, they decompose the weights $W \in \mathbb{R}^{s \times s}$ of a fully-connected layer, acting on $x \in \mathbb{R}^s$ as $\text{vec}(W) = U^G v$ where $U^G \in \mathbb{R}^{s \times s}$ are the “symmetry matrices” and $v \in \mathbb{R}^s$ are the “filter weights”. Then they use meta-learning to learn U^G and during the main training keep U^G fixed and only learn v . We may compare MSR to our approach by setting $d = s$. First, note that although the dimensionality of $U \in \mathbb{R}^{nd \times d}$ seems similar to our $L \in \mathbb{R}^{n \times d \times d}$, the L_i are n matrices of shape $d \times d$, whereas U has shape $(nd) \times d$ with many more parameters than L . Also, the weights of L-conv $W \in \mathbb{R}^{n \times m_l \times m_l}$, with m_l being the number of channels, are generally much fewer than MSR filters $v \in \mathbb{R}^d$. Finally, the way in which Uv acts on data is different from L-conv, as the dimensions reveal. The prohibitively high dimensionality of U requires MSR to adopt a sparse-coding scheme, mainly Kronecker decomposition. Though not necessary, we too choose to use a sparse format for L_i , finding that very low-rank L_i often perform best. A Kronecker decomposition may bias the structure of U^G as it introduces a block structure into it.

Augerino In a concurrent work, (Benton et al., 2020) propose Augerino, a method to learn equivariance with neural networks, but restricted to a subgroup of the augmentation transformations. Augerino learns which subset of the augmentations improved the prediction. This is done by writing The data augmentation is written as $g_\epsilon = \exp(\sum_i \epsilon_i \theta_i L_i)$ (equation (9) in Benton et al. (2020)), with randomly sampled $\epsilon_i \in [-1, 1]$. θ_i are trainable weights which determine which L_i helped with the learning task. Furthermore, Their Lie algebra is fixed to affine transformations in 2D (translations, rotations, scaling and shearing). Our approach is more general. We learn the L_i directly without restricting to known symmetries. Additionally, we do not use the exponential map or matrix logarithm, hence, our method is easy to implement. Lastly, Augerino uses sampling to effectively cover the space of group transformations. Since we work with the Lie algebra rather the group itself, we do not require sampling.

C Experiments

We conduct a set of experiments to see how well L-conv can extract infinitesimal generators.

Nonlinear activation As noted in Weiler et al. (2018a), an arbitrary nonlinear activation σ may not keep the architecture equivariant under G . However, as we showed in SI A (**Extended equivariance for L-conv**), the feature dimensions can pass through any nonlinear neural network without affecting the equivariance of L-conv. This means that the weights of the nonlinear layer should act only on F and not S .

Implementation The basic way to implement L-conv is as multiple parallel GCN units with aggregation function $f(A)$ being (propagation rule) being \hat{L}_i . We do not use Deep Graph Library (DGL) or other libraries, as we want to make \hat{L}_i learnable for discovering symmetries. A more detailed way to implement L-conv is to encode \hat{L}_i in the form of equation 72 to ensure that there is an underlying shared generator $\hat{\ell}_i$ for all μ which is pushed forward using g_μ , shared for all i . To implement L-conv this way, we need the lift g_μ and the edge weights $w_i^\alpha = [\hat{\ell}_i]_\alpha^\alpha$. With known symmetries, the learnable parameters are W^0 and $\bar{\epsilon}^i$. When the topology and hence B is known, we can encode it into the geometry and only learn the edge weights $[\hat{\ell}_i]_\mu^\alpha$, similar to edge features in message passing neural networks (MPNN) (Gilmer et al., 2017). In general each $\hat{\ell}_i$ has $|E_j|$ (i.e. number of edges) components. We can further reduce these using equation 72, where instead of n matrices $\hat{\ell}_i$, we learn one g_μ shared for all, and a small set of elements $[\hat{\ell}_i]_\alpha^\alpha$. This is easiest when the graph is a regular lattice and each vertex has the same number of neighbors. When the topology of the underlying space is not known (e.g. point cloud or scrambled coordinates), we can learn \hat{L}_i as $d \times d$ matrices. We do this for the scrambled image tests, where we encode \hat{L}_i as low-rank matrices.

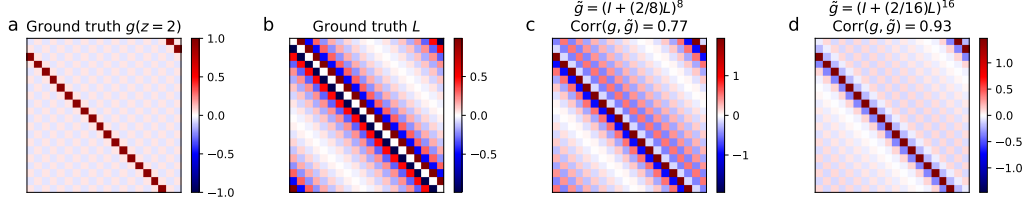


Figure 6: 1D Translation: Using the Shannon-Whittaker Interpolation (SWI) one can generate continuous shifts on discrete data. These include integer shifts (a, ground truth). (SWI) also yields an infinitesimal generator L for shifts (b). This L can be used to approximate finite shifts using $\tilde{g}_n(z) = (I + z/nL)^n$, with $n \rightarrow \infty$ yielding $\exp[zL]$. (c) and (d) show the approximation of a shift by two pixels using $n = 8$ and $n = 16$.

Symmetry Discovery Literature In addition to simplifying the construction of equivariant architectures, our method can also learn the symmetry generators from data. Learning symmetries from data has been studied before, but mostly in restricted settings. Examples include commutative Lie groups as in Cohen & Welling (2014), 2D rotations and translations in Rao & Ruderman (1999), Sohl-Dickstein et al. (2010) or permutations (Anselmi et al., 2019). Zhou et al. (2020) uses meta-learning to automatically learn symmetries in the data. Yet their weight-sharing scheme and the encoding of the symmetry generators is very different from ours. (Benton et al., 2020) propose Augerino, a method to learn equivariance with neural networks, but restricted to a subgroup of the augmentation transformations. Their Lie algebra is fixed to affine transformations in 2D (translations, rotations, scaling and shearing). Our approach is more general. We learn the L_i directly without restricting to known symmetries. Additionally, we do not use the exponential map or matrix logarithm, hence, our method is easy to implement. Lastly, Augerino uses sampling to effectively cover the space of group transformations. Since we work with the Lie algebra rather the group itself, we do not require sampling.

C.1 Approximating 1D CNN

As discussed in the text, Rao & Ruderman (1999, sec. 4) used the Shannon-Whittaker Interpolation (SWI) (Whittaker, 1915) to define continuous translation on periodic 1D arrays as $\mathbf{F}_\rho^d = g(z)_\rho^\nu \mathbf{F}_\nu$. Here $g(z)_\rho^\nu = \frac{1}{d} \prod_{p=d/2}^{d/2} \cos \frac{2\pi p}{d}(z + \rho - \nu)$ approximates the shift operator for continuous z . These $g(z)$ form a 1D translation group G as $g(w)g(z) = g(w+z)$ with $g(0)_\rho^\nu = \delta_\rho^\nu$. For any $z = \mu \in \mathbb{Z}$, $g_\mu = g(z = \mu)$ are circulant matrices that shift by μ as $[g_\mu]_\nu^\rho = \delta_\nu^{\rho - \mu}$. g_μ can be approximated using the Lie algebra and written as multi-layer L-conv as in sec. 3.1. Using $g(0)_\rho^\nu \hat{L}(\rho - \nu)$, the single Lie algebra basis $[\hat{L}]_0 = \partial_z g(z)|_{z=0}$, acts as $\hat{L}f(\mathbf{z}) = \partial_z f(z)$ (because $\partial_z \delta(z - \nu)f(z) = \partial_\nu f(\nu)$). Its components are $\hat{L}_\rho^\nu = L(\rho - \nu) = \prod_p \frac{2\pi p}{d} \sin \frac{2\pi p}{d}(\rho - \nu)$, which are also circulant due to the $(\rho - \nu)$ dependence. Hence, $[\hat{L}\mathbf{F}]_\rho = \sum_\nu L(\rho - \nu) \mathbf{F}_\nu = [L \star \mathbf{F}]_\nu$ is a convolution. Rao & Ruderman (1999) already showed that this \hat{L} can reproduce finite discrete shifts g_μ used in CNN. They used a primitive version of L-conv with $g_\mu = (I + \epsilon \hat{L})^N$. Thus, L-conv can approximate 1D CNN. This result generalizes easily to higher dimensions.

Figure 6 shows how this approximation works. (b) shows the analytical form of \hat{L} . (c) and (d) show two approximations of $g_2 = g(z = 2)$, shift by two pixels, using $\tilde{g}(z) = (I + z/n\hat{L})^n$ with $n = 8$ and $n = 16$. We can evaluate the quality of these approximations using their cosine correlation defined as $\text{Corr}(g, \tilde{g}) = \frac{\text{Tr}(g^T \tilde{g})}{\sqrt{\|g\|_F \|\tilde{g}\|_F}}$ where $\|A\|_F = \sqrt{\sum_{i,j} (A_{ij}^i)^2}$ is the Frobenius L_2 norm. (c) shows 0.77 correlation and (d) has 0.93.

C.2 Extracting 2D rotation generator for fixed small rotations

Ground truth We can use the same SWI 1D translation generators discussed above for CNN as ∂_x and ∂_y to construct the rotation generator $L_\theta = x\partial_y - y\partial_x$ (Fig. 7, a). We will use cosine correlation with this L_θ to evaluate the quality of the learned \hat{L} . As we will find below, the best outcome is from

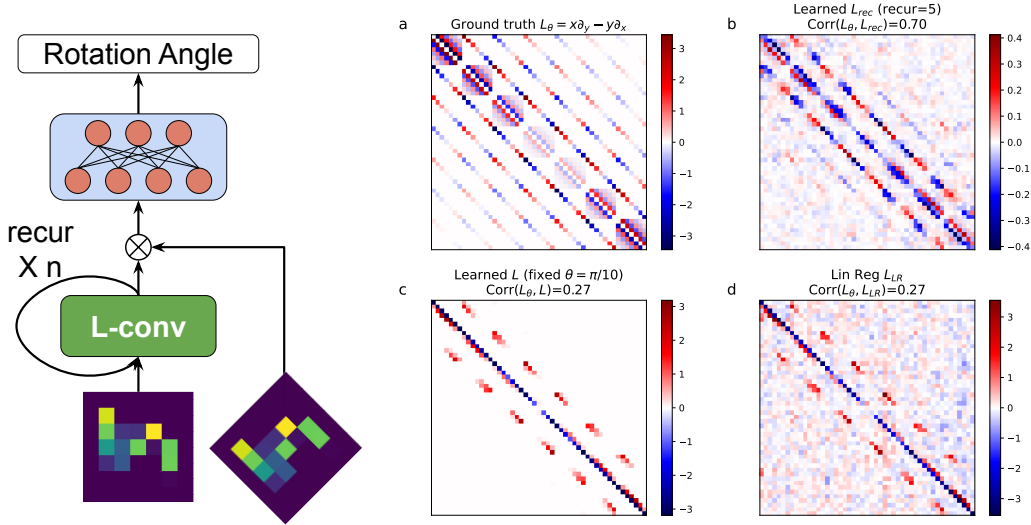


Figure 7: **Learning the infinitesimal generator of $SO(2)$** Left shows the architecture for learning rotation angles between pairs of images. (a) shows the rotation generator calculated analytically using Shannon-Whittaker interpolation. (b) is the \hat{L} learned using recursive L-conv learning rotation angle between a pair of images. (c) is an L learned using a fixed small rotation angle $\theta = \pi/10$, and (d) shows \hat{L} found using the numeric linear regression solution from the fixed angle data. (b) has the highest cosine correlation (0.70) with the ground truth, compared to 0.27 for \hat{L} extracted using small angles.

\hat{L} learned using a recursive L-conv learning the angle of rotation between a pair of images (Fig. 7, b) with 0.70 correlation.

Using fixed small angle In the first experiment we try to learn a small rotation with angle $\theta = \pi/10$ using a single layer L-conv (Fig. 7,c). This experiment was already done in (Rao & Ruderman, 1999). The input is a random 7×7 image \mathbf{f} with pixels chosen in $[-.5, 0.5)$. The output \mathbf{f}^θ is the same image rotated by θ using pytorch affine transform. Our training set contains 50,000 images., the test set was 10,000 images, batch size was 64. The code was implemented in pytorch and we used the Adam optimizer with learning rate 10^{-2} . The experiments were run for 20 epochs. This problem is simply a linear regression with $\mathbf{f}^\theta = R\mathbf{f} = (I + \epsilon\hat{L})\mathbf{f}$. L-conv solves it using SGD and finds $\epsilon\hat{L}$. This problem can also be solved exactly using the solution to linear regression. Let $X = (\mathbf{f}_1, \dots, \mathbf{f}_N)$ and $Y = (\mathbf{f}_1^\theta, \dots, \mathbf{f}_N^\theta)$ be the matrix of all inputs and outputs, respectively. The rotation equation is $Y^T = RX^T$. Thus, the rotation matrix is given by $R = (Y^T X)(X^T X)^{-1}$. Figure 7 (c, d) shows the results of this experiment. The L found using L-conv with SGD is much cleaner than the numerical linear regression solution $L_{LR} = (R - I)/\theta$. The loss becomes extremely small both on training and test data.

C.3 Learning rotation angle

In this experiment we have a pair of input images $(\mathbf{f}_n, R(\theta_n)\mathbf{f}_n)$, with $R(\theta) \in SO(2)$ (approximating 2D rotations). The two inputs differ by a finite rotation with angle $\theta_n \in [0, \pi/8)$. The task is to learn the rotation angle θ_n . For this task we use a recursive L-conv. We set $W^0 = I$. The L-conv weight $\bar{\epsilon}$ is $m \times m$. To be able to encode multiple angles, we set $m = 10$ and feed 10 copies of the \mathbf{f}_n as input $\mathbf{h}_0 = [\mathbf{f}_n]_{10}$. We pass this through the same L-conv layer $t = 3$ times as $\mathbf{h}_i = Q[\mathbf{h}_{i-1}]$. In the final layer, we first take the dot product of the final output \mathbf{h}_t with the rotated input $\mathbf{y}_n = R(\theta)\mathbf{f}_n$ to obtain $\mathbf{g} = \tanh(\mathbf{y}_n^T \mathbf{h}_t)$. We then pass the output $\mathbf{g} \in \mathbb{R}^m$ through a fully-connected (FC) layer with 5 nodes and tanh activation, and finally through a linear FC layer with one output to obtain the angle. The batch size was 16, Adam optimizer, learning rate 10^{-3} , rest were default.

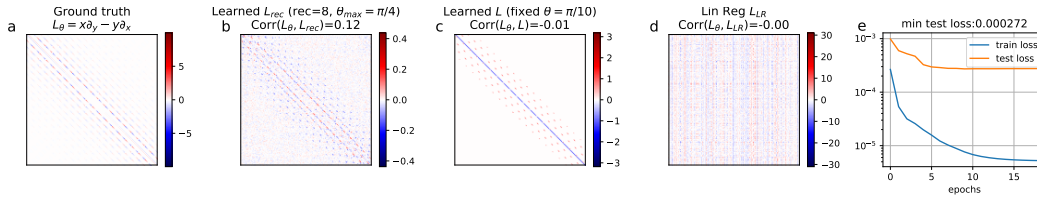


Figure 8: Learning \hat{L} via larger rotation angles for larger images. This time the correlation with ground truth L_θ is much less, but accuracy is still very good.

Despite being a much harder task than fitting a fixed angle rotation, the learned \hat{L} of this experiment has the highest (0.70) cosine correlation with the ground truth L_θ (Fig. 7, b). Even though the architecture is rather complicated and L-conv is followed by two MLP layers, the \hat{L} in L-conv learns the infinitesimal generator of rotations very well. We also conducted experiments with larger random images (20 20) and larger angles of rotation $\theta_n \in [0, \pi/4]$ (Fig. 8). While the accuracy of learning the angles is still pretty good (Fig. 8, e, test loss $2.7 \cdot 10^{-4}$) the larger angles result in less correlation between the learned \hat{L} and the ground truth L_θ (Fig. 8, b, correlation 0.12 with 8 times recurrence). The learned \hat{L} is closer to a finite angle rotation. This may be because the with small number of recurrences the network found small but finite rotations approximate larger rotations better than using a true infinitesimal generator.

D Experiments on Images

To understand precisely how L-conv performs in comparison with CNN and other baselines, we conduct a set of carefully designed experiments. Defining pooling for L-conv merits more research. Without pooling, we cannot use L-conv in state-of-the-art models for problems such as image classification. Therefore, we use the simplest possible models in our experiments: one or two L-conv, or CNN, or FC layers, followed by a classification layer. We do not use any other operations such as dropout or batch normalization in any of the experiments.

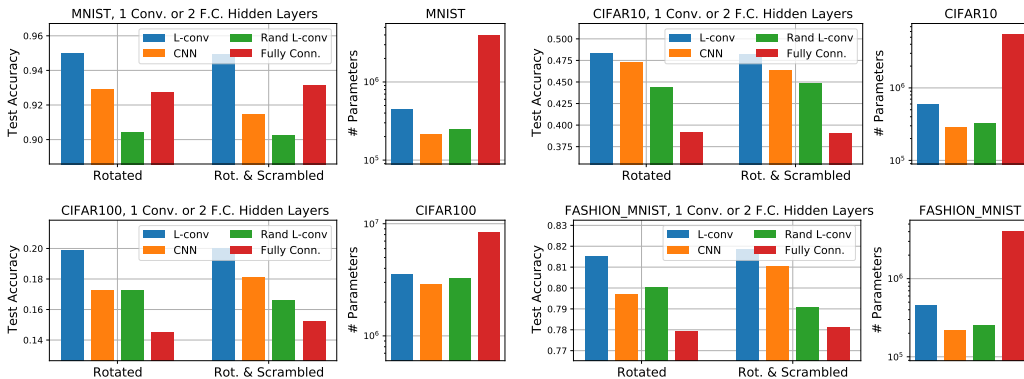


Figure 9: Results on four datasets with two variant: “Rotated” and “Rotated and scrambled”. In all cases L-conv performs best. On MNIST, FC and CNN come close, but using 5x more parameters.

Test Datasets We use four datasets: MNIST, CIFAR10, CIFAR100, and FashionMNIST. To test the efficiency of L-conv in dealing with hidden or unfamiliar symmetries, we conducted our tests on two modified versions of each dataset: 1) **Rotated**: each image rotated by a random angle (no augmentation); 2) **Rotated and Scrambled**: random rotations are followed by a fixed random permutation (same for all images) of pixels. We used a 80-20 training test split on 60,000 MNIST and FashionMNIST, and on 50,000 CIFAR10 and CIFAR100 images. Scrambling destroys the correlations existing between values of neighboring pixels, removing the locality of features in images. As a result, CNN need to encode more patterns, as each image patch has a different correlation pattern.

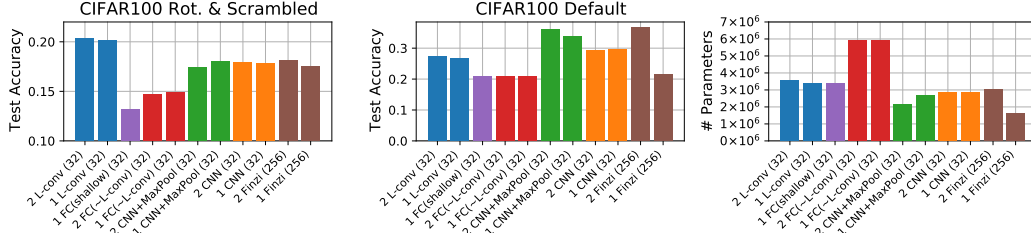


Figure 10: Comparison of one and two layer performance of L-conv (blue), CNN without pooling (orange), CNN with Maxpooling after each layer (green), fully connected (FC) with structure similar to L-conv (red), Lie-Conv (Finzi) Finzi et al. (2020) (brown), and shallow FC, which has a single hidden layer with width such that the total number of parameters matches L-conv (purple). The labels indicate number of layers, layer architecture and number of filters (e.g. “2 L-conv (32)” means two layers of L-conv with 32 filters followed by one classification layer). Left and middle plots show test accuracies on CIFAR100 with rotated and scrambled images, and on the original CIFAR100 dataset, respectively. The plot on the right shows the number of parameters in each model, which is the same for the two datasets.

Test Model Architectures We conduct controlled experiments, with one (Fig. 9) or two (Fig. 10) hidden layers being either L-conv or a baseline, followed by a classification layer. For CNN, L-conv and L-conv with random L_i , we used $n_f = m_l = 32$ for number of output filters (i.e. output dimension of W^i). For CNN we used 3×3 kernels and equivalently used $n_L = 9$ for the number of L_i in L-conv and random L-conv. We also used “LieConv” Finzi et al. (2020) as a baseline (Fig. 10, brown). We used the default $k = 256$ in LieConv, which yields comparable number of parameters to our other models. For the symmetry group in LieConv we used $SE(3)$. We also used the default ResNet architecture provided by Finzi et al. (2020) for both the one and two layer experiments. We turned off batch normalization, consistent with other experiments. We encode L_i as sparse matrices $L_i = U_i V_i$ with hidden dimension $d_h = 16$ in Fig. 9 and $d_h = 8$ in Fig. 10, showing that very sparse L_i can perform well. The weights W^i are each $m_l \times m_{l+1}$ dimensional. The output of the L-conv layer is $d \times m_{l+1}$. As mentioned above, we use two FC baselines. The FC in Fig. 9 and FC (L-conv) in Fig. 10 mimic L-conv, but lacks weight-sharing. The FC weights are $W = ZV$ with V being $(n_L d_h) \times d$ and Z being $(m_{l+1} \times d) \times d_h$. For “FC (shallow)” in Fig. 10, we have one wide hidden layer with $u = n_L \text{ conv} / (m d c)$, where $n_L \text{ conv}$ is the total number of parameters in the L-conv model, m and c the input and output channels, and d is the input dimension. We experimented with encoding L_i as multi-layer perceptrons, but found that a single hidden layer with linear activation works best. We also conduct tests with two layers of L-conv, CNN and FC (Fig. 10), with each L-conv, CNN and FC layer as described above, except that we reduced the hidden dimension in L_i to $d_h = 8$.

Baselines We compare L-conv against four baselines: CNN, random L_i , fully connected (FC) and LieConv. Using CNN or $SE(3)$ LieConv on scrambled images amounts to using poor inductive bias in designing the architecture. Similarly, random, untrained L_i is like using bad inductive biases. Testing on random L_i serves to verify that L-conv’s performance is not due to the structure of the architecture, and that the L_i in L-conv really learn patterns in the data. Finally, to verify that the higher parameter count in L-conv is not responsible for the high performance, we construct two kinds of FC models. The first type (“Fully Conn.” in Fig. 9 and “FC (L-conv)” in Fig. 10) is a multilayer FC network with the same input ($d \times m_0$), hidden ($k \times n_L$ for low-rank L_i) and output ($d \times m_1$) dimensions as L-conv, but lacking the weight-sharing, leading to much larger number parameters than L-conv. The second type (“FC (shallow)” in Fig. 10) consists of a single hidden layer with a width such that the total number of model parameters match L-conv.

Results Fig. 9 shows the results for single layer experiments. On all four datasets both in the rotated and the rotated and scrambled case L-conv performed considerably better than CNN and the baselines. Compared to CNN, L-conv naturally requires extra parameters to encode L_i , but low-rank encoding with rank $d_h \times d$ only requires $O(d_h d)$ parameters, which can be negligible compared to FC layers. We observe that FC layers consistently perform worse than L-conv, despite having much more parameters than L-conv. We also find that not training the L_i (“Rand L-conv”) leads to significant performance drop. We ran tests on the unmodified images as well (Supp. Fig 12), where CNN performed best, but L-conv trails closely behind CNN.

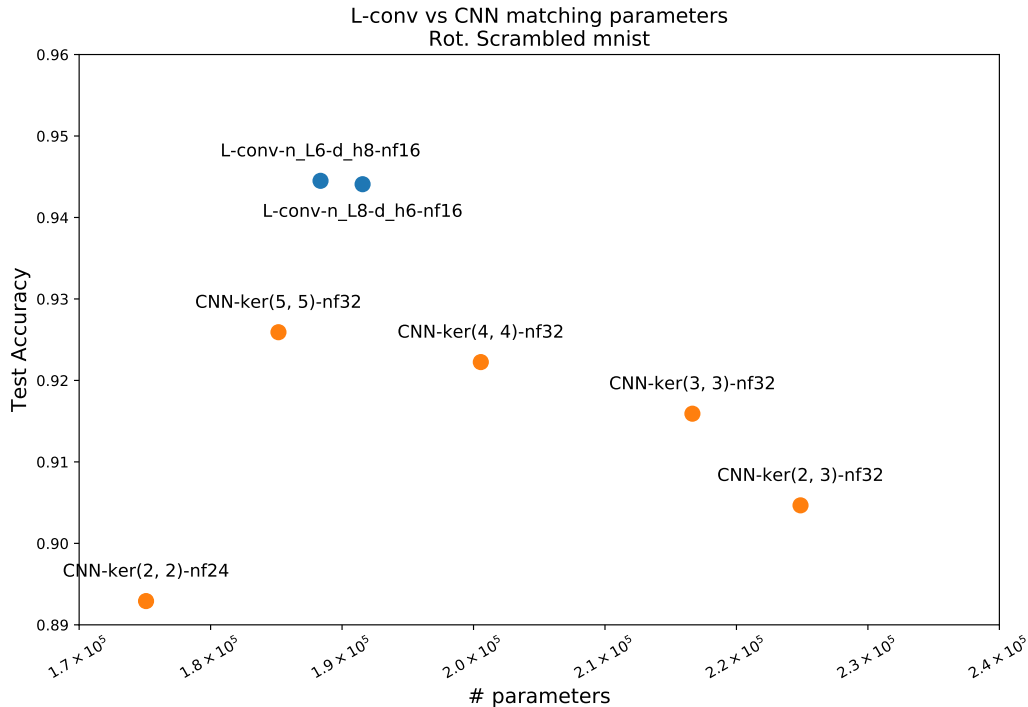


Figure 11: Matching number of parameters in CNN and L-conv, we observe that L-conv still performs better on Rotated and Scrambled MNIST.

Additional experiments testing the effect of number of layers, number of parameters and pooling are shown in Fig. 10. On CIFAR100, we find that both FC configurations, FC(L-conv) and FC(shallow) consistently perform worse than L-conv, evidence that L-conv’s performance is *not* due to its extra parameters. L-conv outperforms all other tested models on rotated and scrambled CIFAR100, including LieConv. Without pooling, we observe that both L-conv and CNN do not benefit from adding a second layer. On the default CIFAR100 dataset, one and two layer CNN with max-pooling perform significantly better than L-Conv. Two Layer $SE(3)$ LieConv (labelled “2 Finzi (256)”) performs best on default CIFAR100, but not on the scrambled and rotated version. This is expected, as the symmetries of the latter are masked by the scrambling. This is where the benefit of our model becomes evident, namely cases where the data may have hidden or unfamiliar symmetries. We also verified that the higher performance of L-conv compared to CNN is not due to higher number of parameters (Appendix D.2)

D.1 Details of experiments

Hardware and Implementation We implemented L-conv in Keras and Tensorflow 2.2 and ran our tests on a system with a 6 core Intel Core i7 CPU, 32GB RAM, and NVIDIA Quadro P6000 (24GB RAM) GPU. The L-conv layer did not require significantly more resources than CNN and ran only slightly slower.

D.2 Additional Experiments

Matching number of parameters in CNN To verify that the difference in the number of parameters between CNN and L-conv was not responsible for the improved performance, we ran experiment where we allowed the kernel-size of L-conv and CNN to differ and tried to match the number of parameters between the two. Fig. 11 shows that on rotated and scrambled MNIST L-conv still performs better than CNN even after the latter has been allowed to have the same or more number of parameters than L-conv.

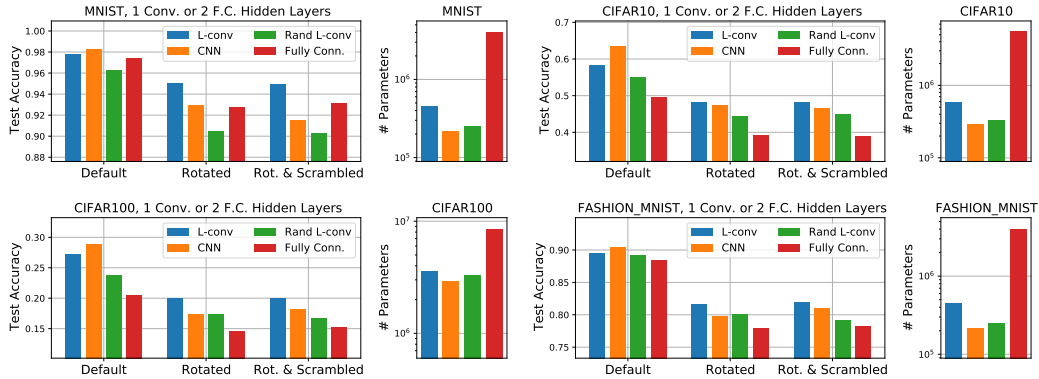


Figure 12: Test results on four datasets with three variant: “Default” (unmodified dataset), “Rotated” and “Rotated and scrambled”. On the Default dataset, CNN performs best, but L-conv is always the second best. For Rotated and Rot. & Scrambled, in all cases L-conv performed best. In MNIST, FC and CNN layers come close, but using 5x more parameters.

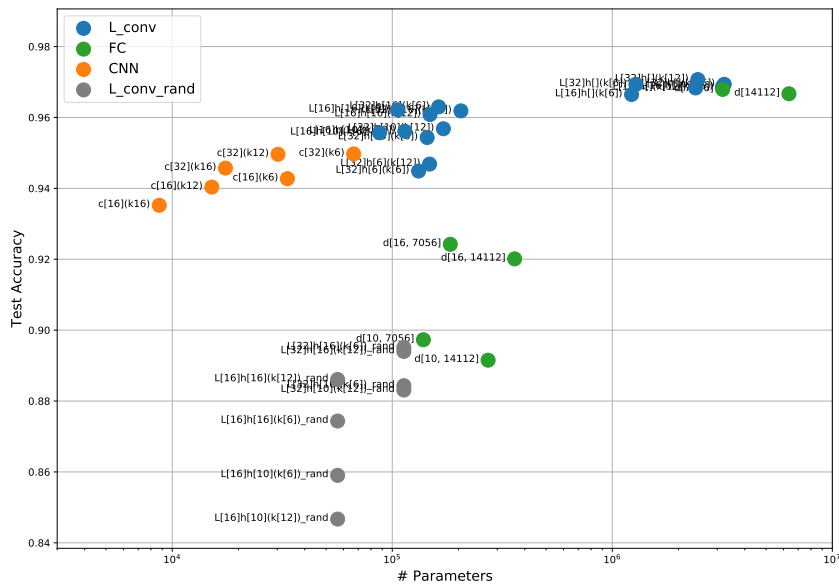


Figure 13: Training low-rank L-conv layer during training.

In Figure 13 we compare the performance of a single layer of L-conv on a classification task on scrambled rotated MNIST, where pixels have been permuted randomly and images have been rotated between -90 to $+90$ degrees. The models consisted of a final classification layer preceded by either one L-conv (blue), or one CNN (orange), or multiple fully-connected (FC, green) layers with similar number of neurons as the L-conv, but without weight sharing. We see that most L-conv configurations had the highest performance without a too many trainable parameters. Note that, parameters in FC layers are much higher than comparable L-conv, but yield worse results. The dots are labeled to show the configurations, with $L[32]h[6](k[6])$ meaning $k = 6$ as number of L_i , 32 output filters, and $h = 6$ hidden dimensions for low-rank encoding of L_i . The y-axis shows the test accuracy and the x-axis the number of trainable parameters. The grey lines show the performance of L-conv with fixed random L_i , but trainable shared weights, showing that indeed the learned L_i improve the performance quite significantly.