



# Universal Approximations of Invariant Maps by Neural Networks

Dmitry Yarotsky<sup>1,2</sup>

Received: 15 June 2019 / Accepted: 20 August 2020 / Published online: 30 April 2021  
© The Author(s), under exclusive licence to Springer Science+Business Media, LLC, part of Springer Nature 2021

## Abstract

We describe generalizations of the universal approximation theorem for neural networks to maps invariant or equivariant with respect to linear representations of groups. Our goal is to establish network-like computational models that are both invariant/equivariant and provably complete in the sense of their ability to approximate any continuous invariant/equivariant map. Our contribution is three-fold. First, in the general case of compact groups we propose a construction of a complete invariant/equivariant network using an intermediate polynomial layer. We invoke classical theorems of Hilbert and Weyl to justify and simplify this construction; in particular, we describe an explicit complete ansatz for approximation of permutation-invariant maps. Second, we consider groups of translations and prove several versions of the universal approximation theorem for convolutional networks in the limit of continuous signals on euclidean spaces. Finally, we consider 2D signal transformations equivariant with respect to the group  $SE(2)$  of rigid euclidean motions. In this case we introduce the “charge-conserving convnet”—a convnet-like computational model based on the decomposition of the feature space into isotypic representations of  $SO(2)$ . We prove this model to be a universal approximator for continuous  $SE(2)$ —equivariant signal transformations.

**Keywords** Neural network · Approximation · Linear representation · Invariance · Equivariance · Polynomial · Polarization · Convnet

**Mathematics Subject Classification** 41A63 · 41A29 · 68Q17 · 62J02 · 62M45 · 20C35 · 13A50

---

Communicated by Wolfgang Dahmen, Ronald A. Devore, and Philipp Grohs.

---

✉ Dmitry Yarotsky  
d.yarotsky@skoltech.ru

<sup>1</sup> Skolkovo Institute of Science and Technology, Moscow, Russia

<sup>2</sup> Institute for Information Transmission Problems, Moscow, Russia

# 1 Introduction

## 1.1 Motivation

**Symmetric models** An important topic in learning theory is the design of predictive models properly reflecting symmetries naturally present in the data (see, e.g., [3, 35, 37]). Most commonly, in the standard context of supervised learning, this means that our predictive model should be *invariant* with respect to a suitable *group of transformations*: given an input object, we often know that its class or some other property that we are predicting does not depend on the object representation (e.g., associated with a particular coordinate system), or for other reasons does not change under certain transformations. In this case we would naturally like the predictive model to reflect this independence. If  $f$  is our predictive model and  $\Gamma$  the group of transformations, we can express the property of invariance by the identity  $f(\mathcal{A}_\gamma \mathbf{x}) = f(\mathbf{x})$ , where  $\mathcal{A}_\gamma \mathbf{x}$  denotes the action of the transformation  $\gamma \in \Gamma$  on the object  $\mathbf{x}$ .

There is also a more general scenario where the output of  $f$  is another complex object that is supposed to transform appropriately if the input object is transformed. This scenario is especially relevant in the setting of multi-layered (or stacked) predictive models, if we want to propagate the symmetry through the layers. In this case one speaks about *equivariance*, and mathematically it is described by the identity  $f(\mathcal{A}_\gamma \mathbf{x}) = \mathcal{A}_\gamma f(\mathbf{x})$ , assuming that the transformation  $\gamma$  acts in some way not only on inputs, but also on outputs of  $f$ . (For brevity, here and in the sequel we will slightly abuse notation and denote any action of  $\gamma$  by  $\mathcal{A}_\gamma$ , though of course in general the input and output objects are different and  $\gamma$  acts differently on them. It will be clear which action is meant in a particular context).

A well-known important example of equivariant transformations are convolutional layers in neural networks, where the group  $\Gamma$  is the group of grid translations,  $\mathbb{Z}^d$ .

**Symmetrization vs. intrinsic symmetry** We find it convenient to roughly distinguish two conceptually different approaches to the construction of invariant and equivariant models that we refer to as the *symmetrization-based* one and the *intrinsic* one. The symmetrization-based approach consists in starting from some asymmetric model, and symmetrizing it by a group averaging. On the other hand, the intrinsic approach consists in imposing prior structural constraints on the model that guarantee its symmetry.

In the general mathematical context, the difference between the two approaches is best illustrated with the example of *symmetric polynomials* in the variables  $x_1, \dots, x_n$ , i.e., the polynomials invariant with respect to arbitrary permutations of these variables. With the symmetrization-based approach, we can obtain any invariant polynomial by starting with an arbitrary polynomial  $f$  and symmetrizing it over the group of permutations  $S_n$ , i.e. by defining  $f_{\text{sym}}(x_1, \dots, x_n) = \frac{1}{n!} \sum_{\rho \in S_n} f(x_{\rho(1)}, \dots, x_{\rho(n)})$ . On the other hand, the intrinsic approach is associated with the fundamental theorem of symmetric polynomials, which states that any invariant polynomial  $f_{\text{sym}}$  in  $n$  variables can be obtained as a superposition  $f(s_1, \dots, s_n)$  of some polynomial  $f$  and the *elementary symmetric polynomials*  $s_1, \dots, s_n$ . Though both approaches yield essentially

the same result (an arbitrary symmetric polynomial), the two constructions are clearly very different.

In practical machine learning, symmetrization is ubiquitous. It is often applied both on the level of data and the level of models. This means that, first, prior to learning an invariant model, one augments the available set of training examples  $(\mathbf{x}, f(\mathbf{x}))$  by new examples of the form  $(\mathcal{A}_\gamma \mathbf{x}, f(\mathbf{x}))$  (see, for example, Section B.2 of Thoma [43] for a list of transformations routinely used to augment datasets for image classification problems). Second, once some, generally non-symmetric, predictive model  $\hat{f}$  has been learned, it is symmetrized by setting  $\hat{f}_{\text{sym}}(\mathbf{x}) = \frac{1}{|\Gamma_0|} \sum_{\gamma \in \Gamma_0} \hat{f}(\mathcal{A}_\gamma \mathbf{x})$ , where  $\Gamma_0$  is some subset of  $\Gamma$  (e.g., randomly sampled). This can be seen as a manifestation of the symmetrization-based approach, and its practicality probably stems from the fact that the real world symmetries are usually only approximate, and in this approach one can easily account for their imperfections (e.g., by adjusting the subset  $\Gamma_0$ ). On the other hand, the weight sharing in convolutional networks [22,45] can be seen as a manifestation of the intrinsic approach (since the translational symmetry is built into the architecture of the network from the outset), and convnets are ubiquitous in modern machine learning [23].

**Completeness** In this paper we will be interested in the theoretical opportunities of the intrinsic approach in the context of approximations using neural-network-type models. Suppose, for example, that  $f$  is an invariant map that we want to approximate with the usual ansatz of a perceptron with a single hidden layer,  $\hat{f}(x_1, \dots, x_d) = \sum_{n=1}^N c_n \sigma(\sum_{k=1}^d w_{nk} x_k + h_n)$  with some nonlinear activation function  $\sigma$ . Obviously, this ansatz breaks the symmetry, in general. Our goal is to modify this ansatz in such a way that, first, it does not break the symmetry and, second, it is *complete* in the sense that it is not too specialized and any reasonable invariant map can be arbitrarily well approximated by it. In Sect. 2 we show how this can be done by introducing an extra polynomial layer into the model. In Sects. 3, 4 we will consider more complex, deep models (convnets and their modifications). We will understand completeness in the sense of the universal approximation theorem for neural networks [32].

**Linear representations** Designing invariant and equivariant models requires us to decide how the symmetry information is encoded in the layers. A standard assumption, to which we also will adhere in this paper, is that the group acts by linear transformations. Precisely, when discussing invariant models we are looking for maps of the form

$$f : V \rightarrow \mathbb{R}, \tag{1.1}$$

where  $V$  is a vector space carrying a linear representation  $R : \Gamma \rightarrow \text{GL}(V)$  of a group  $\Gamma$ . More generally, in the context of multi-layer models

$$f : V_1 \xrightarrow{f_1} V_2 \xrightarrow{f_2} \dots \tag{1.2}$$

we assume that the vector spaces  $V_k$  carry linear representations  $R_k : \Gamma \rightarrow \text{GL}(V_k)$  (the “baseline architecture” of the model), and we must then ensure equivariance in

each link. Note that a linear action of a group on the input space  $V_1$  is a natural and general phenomenon. In particular, the action is linear if  $V_1$  is a linear space of functions on some domain, and the action is induced by (not necessarily linear) transformations of the domain. Prescribing linear representations  $R_k$  is then a viable strategy to encode and upkeep the symmetry in subsequent layers of the model.

**Compact vs. non-compact groups** From the perspective of approximation theory, we will be interested in *finite* computational models, i.e. including finitely many operations as performed on a standard computer. Finiteness is important for potential studies of approximation rates (though such a study is not attempted in the present paper). Compact groups have the nice property that their irreducible linear representations are finite-dimensional. This allows us, in the case of such groups, to modify the standard shallow neural network ansatz so as to obtain a computational model that is finite, fully invariant/equivariant and complete, see Sect. 2. On the other hand, irreducible representations of non-compact groups such as  $\mathbb{R}^\nu$  are infinite-dimensional in general. As a result, finite computational models can be only approximately  $\mathbb{R}^\nu$ -invariant/equivariant. Nevertheless, we show in Sects. 3, 4 that complete  $\mathbb{R}^\nu$ —and  $SE(\nu)$ —equivariant models can be rigorously described in terms of appropriate limits of finite models.

## 1.2 Related Work

Our work can be seen as an extension of results on the universal approximation property of neural networks [7,10,18,19,24,29,31,32] to the setting of group invariant/equivariant maps and/or infinite-dimensional input spaces.

Our general results in Sect. 2 are based on classical results of the theory of polynomial invariants [16,17,46].

An important element of constructing invariant and equivariant models is the extraction of invariant and equivariant features. In the present paper we do not focus on this topic, but it has been studied extensively, see e.g. general results along with applications to 2D and 3D pattern recognition in [3,27,35,37,41].

In a series of works reviewed in Cohen et al. [5], the authors study expressiveness of deep convolutional networks using hierarchical tensor decompositions and convolutional arithmetic circuits. In particular, representation universality of several network structures is examined in Cohen and Shashua [4].

In a series of works reviewed in Poggio et al. [33], the authors study expressiveness of deep networks from the perspective of approximation theory and hierarchical decompositions of functions. Learning of invariant data representations and its relation to information processing in the visual cortex has been discussed in Anselmi et al. [1].

In the series of papers [2,25,26,39], multiscale wavelet-based group invariant scattering operators and their applications to image recognition have been studied.

There is a large body of work proposing specific constructions of networks for applied group invariant recognition problems, in particular image recognition approximately invariant with respect to the group of rotations or some of its subgroups: deep symmetry networks of Gens and Domingos [11], G-CNNs of Cohen and Welling [6],

networks with extra slicing operations in Dieleman et al. [8], RotEqNets of Marcos et al. [28], networks with warped convolutions in Henriques and Vedaldi [15], Polar Transformer Networks of Esteves et al. [9]. In Sect. 4 we study a family of models equivariant w.r.t. 2D euclidean motions; our construction partly resembles the one used in Worrall et al. [48]. However, in contrast to all these papers, we are primarily interested in the theoretical guarantees of invariance and completeness.

Our Theorem 3.1 resembles the Curtis-Hedlund-Lyndon theorem from the theory of cellular automata, that states that a map  $f : \{1, \dots, N\}^{\mathbb{Z}^v} \rightarrow \{1, \dots, N\}^{\mathbb{Z}^v}$  is  $\mathbb{Z}^v$ -equivariant and continuous in the product topology if and only if it is defined by a finite cellular automaton [14]. In Theorem 3.1 we characterize the maps  $f : L^2(\mathbb{R}^v, \mathbb{R}^{d_v}) \rightarrow L^2(\mathbb{R}^v, \mathbb{R}^{d_u})$  that are  $\mathbb{R}^v$ -equivariant and continuous in the norm topology as limit points of convnets.

In Sect. 2.4 we apply the invariant theory to construct complete permutation-invariant networks. See [34,49] for related results and discussions of permutation-invariant models, as well as applications to image recognition problems.

In Kondor and Trivedi [20], it is proved that network layers of the conventional structure “a linear transformation followed by pointwise nonlinear activation” are group-equivariant iff the linear part is a (generalized) convolution. This can be viewed as a completeness result for equivariant maps implementable by a single standard layer. In the present paper, our point of view is rather different: first, we strive to describe approximations to maps from general functional classes, and second, we are particularly interested in symmetries like  $\mathbb{R}^v$  or  $\text{SE}(2)$ , that cannot be implemented in a single finite layer, but can be recovered in a suitable limit (cf. Sects. 3, 4).

### 1.3 Contribution of this Paper

As discussed above, we will be interested in the following general question: assuming there is a “ground truth” invariant or equivariant map  $f$ , how can we “intrinsically” approximate it by a neural-network-like model? Our goal is to describe models that are finite, invariant/ equivariant (up to limitations imposed by the finiteness of the model) and provably complete in the sense of approximation theory.

Our contribution is three-fold:

- In Sect. 2 we consider general compact groups and approximations by shallow networks. Using the classical polynomial invariant theory, we describe a general construction of shallow networks with an extra polynomial layer which are exactly invariant/equivariant and complete (Propositions 2.3, 2.4). Then, we discuss how this construction can be improved using the idea of polarization and a theorem of Weyl (Propositions 2.5, 2.7). Finally, as a particular illustration of the “intrinsic” framework, we consider maps invariant with respect to the symmetric group  $S_N$ , and describe a corresponding neural network model which is  $S_N$ -invariant and complete (Theorem 2.4). This last result is based on another theorem of Weyl.
- In Sect. 3 we prove several versions of the universal approximation theorem for convolutional networks and groups of translations. The main novelty of these results is that we approximate maps  $f$  defined on the infinite-dimensional space of *continuous signals* on  $\mathbb{R}^v$ . Specifically, one of these versions (Theorem 3.1)

states that a signal transformation  $f : L^2(\mathbb{R}^v, \mathbb{R}^{d_v}) \rightarrow L^2(\mathbb{R}^v, \mathbb{R}^{d_v})$  can be approximated, in some natural sense, by convnets without pooling if and only if  $f$  is continuous and translationally-equivariant (here, by  $L^2(\mathbb{R}^v, \mathbb{R}^d)$  we denote the space of square-integrable functions  $\Phi : \mathbb{R}^v \rightarrow \mathbb{R}^d$ ). Another version (Theorem 3.2) states that a map  $f : L^2(\mathbb{R}^v, \mathbb{R}^{d_v}) \rightarrow \mathbb{R}$  can be approximated by convnets with pooling if and only if  $f$  is continuous.

- In Sect. 4 we describe a convnet-like model which is a universal approximator for signal transformations  $f : L^2(\mathbb{R}^2, \mathbb{R}^{d_v}) \rightarrow L^2(\mathbb{R}^2, \mathbb{R}^{d_v})$  equivariant with respect to the group  $SE(2)$  of rigid two-dimensional euclidean motions. We call this model *charge-conserving convnet*, based on a 2D quantum mechanical analogy (conservation of the total angular momentum). The crucial element of the construction is that the operation of the network is consistent with the decomposition of the feature space into isotypic representations of  $SO(2)$ . We prove in Theorem 4.1 that a transformation  $f : L^2(\mathbb{R}^2, \mathbb{R}^{d_v}) \rightarrow L^2(\mathbb{R}^2, \mathbb{R}^{d_v})$  can be approximated by charge-conserving convnets if and only if  $f$  is continuous and  $SE(2)$ -equivariant.

## 2 Compact Groups and Shallow Approximations

In this section we give several results on invariant/equivariant approximations by neural networks in the context of compact groups, finite-dimensional representations, and shallow networks. We start by describing the standard group-averaging approach in Sect. 2.1. In Sect. 2.2 we describe an alternative approach, based on the invariant theory. In Sect. 2.3 we show how one can improve this approach using polarization. Finally, in Sect. 2.4 we describe an application of this approach to the symmetric group  $S_N$ .

### 2.1 Approximations Based on Symmetrization

We start by recalling the universal approximation theorem, which will serve as a “template” for our invariant and equivariant analogs. There are several versions of this theorem (see the survey Pinkus [32]), we will use the general and easy-to-state version given in Pinkus [32].

**Theorem 2.1** (Pinkus [32], Theorem 3.1) *Let  $\sigma : \mathbb{R} \rightarrow \mathbb{R}$  be a continuous activation function that is not a polynomial. Let  $V = \mathbb{R}^d$  be a real finite dimensional vector space. Then any continuous map  $f : V \rightarrow \mathbb{R}$  can be approximated, in the sense of uniform convergence on compact sets, by maps  $\hat{f} : V \rightarrow \mathbb{R}$  of the form*

$$\hat{f}(x_1, \dots, x_d) = \sum_{n=1}^N c_n \sigma \left( \sum_{s=1}^d w_{ns} x_s + h_n \right) \quad (2.1)$$

with some coefficients  $c_n, w_{ns}, h_n$ .

Throughout the paper, we assume, as in Theorem 2.1, that  $\sigma : \mathbb{R} \rightarrow \mathbb{R}$  is some (fixed) continuous activation function that is not a polynomial.

Also, as in this theorem, we will understand approximation in the sense of uniform approximation on compact sets, i.e. meaning that for any compact  $K \subset V$  and any  $\epsilon > 0$  one can find an approximating map  $\widehat{f}$  such that  $|f(\mathbf{x}) - \widehat{f}(\mathbf{x})| \leq \epsilon$  (or  $\|f(\mathbf{x}) - \widehat{f}(\mathbf{x})\| \leq \epsilon$  in the case of vector-valued  $f$ ) for all  $\mathbf{x} \in K$ . In the case of finite-dimensional spaces  $V$  considered in the present section, one can equivalently say that there is a sequence of approximating maps  $\widehat{f}_n$  uniformly converging to  $f$  on any compact set. Later, in Sects. 3, 4, we will consider infinite-dimensional signal spaces  $V$  for which such an equivalence does not hold. Nevertheless, we will use the concept of uniform approximation on compact sets as a guiding principle in our precise definitions of approximation in that more complex setting.

Now suppose that the space  $V$  carries a linear representation  $R$  of a group  $\Gamma$ . Assuming  $V$  is finite-dimensional, this means that  $R$  is a homomorphism of  $\Gamma$  to the group of linear automorphisms of  $V$ :

$$R : \Gamma \rightarrow \text{GL}(V).$$

In the present section we will assume that  $\Gamma$  is a *compact* group, meaning, as is customary, that  $\Gamma$  is a compact Hausdorff topological space and the group operations (multiplication and inversion) are continuous. Accordingly, the representation  $R$  is also assumed to be continuous. We remark that an important special case of compact groups are the finite groups (with respect to the discrete topology).

One important property of compact groups is the existence of a unique, both left- and right-invariant Haar measure normalized so that the total measure of  $\Gamma$  equals 1. Another property is that any continuous representation of a compact group on a separable (but possibly infinite-dimensional) Hilbert space can be decomposed into a countable direct sum of irreducible finite-dimensional representations. There are many group representation textbooks to which we refer the reader for details, see e.g. [38,40,44]. Accordingly, in the present section we will restrict ourselves to finite-dimensional representations. Later, in Sects. 3 and 4, we will consider the noncompact groups  $\mathbb{R}^v$  and  $\text{SE}(v)$  and their natural representations on the infinite-dimensional space  $L^2(\mathbb{R}^v)$ , which cannot be decomposed into countably many irreducibles.

Motivated by applications to neural networks, in this section and Sect. 3 we will consider only representations over the field  $\mathbb{R}$  of reals (i.e. with  $V$  a real vector space). Later, in Sect. 4, we will consider complexified spaces as this simplifies the exposition of the invariant theory for the group  $\text{SO}(2)$ .

For brevity, we will call a vector space carrying a linear representation of a group  $\Gamma$  a  $\Gamma$ -*module*. We will denote by  $R_\gamma$  the linear automorphism obtained by applying  $R$  to  $\gamma \in \Gamma$ . In particular, the property that  $R$  is a homomorphism between  $\Gamma$  and  $\text{GL}(V)$  can then be written as

$$R_{\gamma^{-1}} = R_\gamma^{-1}, \quad R_{\gamma\theta} = R_\gamma R_\theta, \quad \forall \gamma, \theta \in \Gamma.$$

The integral over the normalized Haar measure on a compact group  $\Gamma$  is denoted by  $\int_\Gamma \cdot d\gamma$ . We will denote vectors by boldface characters; scalar components of the vector  $\mathbf{x}$  are denoted  $x_k$ .

Recall that given a  $\Gamma$ -module  $V$ , we call a map  $f : V \rightarrow \mathbb{R}$   $\Gamma$ -invariant (or simply invariant) if  $f(R_\gamma \mathbf{x}) = f(\mathbf{x})$  for all  $\gamma \in \Gamma$  and  $\mathbf{x} \in V$ . We state now the basic result on invariant approximation, obtained by symmetrization (group averaging).

**Proposition 2.1** *Let  $\Gamma$  be a compact group and  $V$  a finite-dimensional  $\Gamma$ -module. Then, any continuous invariant map  $f : V \rightarrow \mathbb{R}$  can be approximated by  $\Gamma$ -invariant maps  $\widehat{f} : V \rightarrow \mathbb{R}$  of the form*

$$\widehat{f}(\mathbf{x}) = \int_{\Gamma} \sum_{n=1}^N c_n \sigma(l_n(R_\gamma \mathbf{x}) + h_n) d\gamma, \tag{2.2}$$

where  $c_n, h_n \in \mathbb{R}$  are some coefficients and  $l_n \in V^*$  are some linear functionals on  $V$ , i.e.  $l_n(\mathbf{x}) = \sum_k w_{nk} x_k$ .

**Proof** It is clear that the map (2.2) is  $\Gamma$ -invariant, and we only need to prove the completeness part. Let  $K$  be a compact subset in  $V$ , and  $\epsilon > 0$ . Consider the symmetrization of  $K$  defined by  $K_{\text{sym}} = \cup_{\gamma \in \Gamma} R_\gamma(K)$ . Note that  $K_{\text{sym}}$  is also a compact set, because it is the image of the compact set  $\Gamma \times K$  under the continuous map  $(\gamma, \mathbf{x}) \mapsto R_\gamma \mathbf{x}$ . We can use Theorem 2.1 to find a map  $f_1 : V \rightarrow \mathbb{R}$  of the form  $f_1(\mathbf{x}) = \sum_{n=1}^N c_n \sigma(l_n(\mathbf{x}) + h_n)$  and such that  $|f(\mathbf{x}) - f_1(\mathbf{x})| \leq \epsilon$  on  $K_{\text{sym}}$ . Now consider the  $\Gamma$ -invariant group-averaged map  $\widehat{f}(\mathbf{x}) = \int_{\Gamma} f_1(R_\gamma \mathbf{x}) d\gamma$ . Then for any  $\mathbf{x} \in K$ ,

$$|\widehat{f}(\mathbf{x}) - f(\mathbf{x})| = \left| \int_{\Gamma} (f_1(R_\gamma \mathbf{x}) - f(R_\gamma \mathbf{x})) d\gamma \right| \leq \int_{\Gamma} |f_1(R_\gamma \mathbf{x}) - f(R_\gamma \mathbf{x})| d\gamma \leq \epsilon,$$

where we have used the invariance of  $f$  and the fact that  $|f_1(\mathbf{x}) - f(\mathbf{x})| \leq \epsilon$  for  $\mathbf{x} \in K_{\text{sym}}$ . □

Now we establish a similar result for equivariant maps. Let  $V, U$  be two  $\Gamma$ -modules. For brevity, we will denote by  $R$  the representation of  $\Gamma$  in either of them (it will be clear from the context which one is meant). We call a map  $f : V \rightarrow U$   $\Gamma$ -equivariant if  $f(R_\gamma \mathbf{x}) = R_\gamma f(\mathbf{x})$  for all  $\gamma \in \Gamma$  and  $\mathbf{x} \in V$ .

**Proposition 2.2** *Let  $\Gamma$  be a compact group and  $V$  and  $U$  two finite-dimensional  $\Gamma$ -modules. Then, any continuous  $\Gamma$ -equivariant map  $f : V \rightarrow U$  can be approximated by  $\Gamma$ -equivariant maps  $\widehat{f} : V \rightarrow U$  of the form*

$$\widehat{f}(\mathbf{x}) = \int_{\Gamma} \sum_{n=1}^N R_\gamma^{-1} \mathbf{y}_n \sigma(l_n(R_\gamma \mathbf{x}) + h_n) d\gamma, \tag{2.3}$$

with some coefficients  $h_n \in \mathbb{R}$ , linear functionals  $l_n \in V^*$ , and vectors  $\mathbf{y}_n \in U$ .

**Proof** The proof is analogous to the proof of Proposition 2.1. Fix any norm  $\|\cdot\|$  in  $U$ . Given a compact set  $K$  and  $\epsilon > 0$ , we construct the compact set  $K_{\text{sym}} = \cup_{\gamma \in \Gamma} R_\gamma(K)$  as before. Next, we find  $f_1 : V \rightarrow U$  of the form  $f_1(\mathbf{x}) = \sum_{n=1}^N \mathbf{y}_n \sigma(l_n(\mathbf{x}) + h_n)$



and such that  $\|f(\mathbf{x}) - f_1(\mathbf{x})\| \leq \epsilon$  on  $K_{\text{sym}}$  (we can do it, for example, by considering scalar components of  $f$  with respect to some basis in  $U$ , and approximating these components using Theorem 2.1). Finally, we define the symmetrized map by  $\widehat{f}(\mathbf{x}) = \int_{\Gamma} R_{\gamma}^{-1} f_1(R_{\gamma}\mathbf{x}) d\gamma$ . This map is  $\Gamma$ -equivariant, and, for any  $\mathbf{x} \in K$ ,

$$\begin{aligned} \|\widehat{f}(\mathbf{x}) - f(\mathbf{x})\| &= \left\| \int_{\Gamma} (R_{\gamma}^{-1} f_1(R_{\gamma}\mathbf{x}) - R_{\gamma}^{-1} f(R_{\gamma}\mathbf{x})) d\gamma \right\| \\ &\leq \max_{\gamma \in \Gamma} \|R_{\gamma}\| \int_{\Gamma} \|f_1(R_{\gamma}\mathbf{x}) - f(R_{\gamma}\mathbf{x})\| d\gamma \\ &\leq \epsilon \max_{\gamma \in \Gamma} \|R_{\gamma}\|. \end{aligned}$$

By continuity of  $R$  and compactness of  $\Gamma$ ,  $\max_{\gamma \in \Gamma} \|R_{\gamma}\| < \infty$ , so we can approximate  $f$  by  $\widehat{f}$  on  $K$  with any accuracy. □

Propositions 2.1, 2.2 present the “symmetrization-based” approach to constructing invariant/equivariant approximations relying on the shallow neural network ansatz (2.1). The approximating expressions (2.2), (2.3) are  $\Gamma$ -invariant/equivariant and universal. Moreover, in the case of finite groups the integrals in these expressions are finite sums, i.e. these approximations consist of finitely many arithmetic operations and evaluations of the activation function  $\sigma$ . In the case of infinite groups, the integrals can be approximated by sampling the group.

In the remainder of Sect. 2 we will pursue an alternative approach to symmetrize the neural network ansatz, based on the theory of polynomial invariants.

We finish this subsection with the following general observation. Suppose that we have two  $\Gamma$ -modules  $U, V$ , and  $U$  can be decomposed into  $\Gamma$ -invariant submodules:  $U = \bigoplus_{\beta} U_{\beta}^{m_{\beta}}$  (where  $m_{\beta}$  denotes the multiplicity of  $U_{\beta}$  in  $U$ ). Then a map  $f : V \rightarrow U$  is equivariant if and only if it is equivariant in each component  $U_{\beta}$  of the output space. Moreover, if we denote by  $\text{Equiv}(V, U)$  the space of continuous equivariant maps  $f : V \rightarrow U$ , then

$$\text{Equiv}\left(V, \bigoplus_{\beta} U_{\beta}^{m_{\beta}}\right) = \bigoplus_{\beta} \text{Equiv}(V, U_{\beta})^{m_{\beta}}. \tag{2.4}$$

This shows that the task of describing equivariant maps  $f : V \rightarrow U$  reduces to the task of describing equivariant maps  $f : V \rightarrow U_{\beta}$ . In particular, describing vector-valued invariant maps  $f : V \rightarrow \mathbb{R}^{d_U}$  reduces to describing scalar-valued invariant maps  $f : V \rightarrow \mathbb{R}$ .

### 2.2 Approximations Based on Polynomial Invariants

The invariant theory seeks to describe *polynomial invariants* of group representations, i.e. polynomial maps  $f : V \rightarrow \mathbb{R}$  such that  $f(R_{\gamma}\mathbf{x}) \equiv f(\mathbf{x})$  for all  $\mathbf{x} \in V$ . A fundamental result of the invariant theory is Hilbert’s finiteness theorem [16,17] stating that for completely reducible representations, all the polynomial invariants are alge-

braically generated by a finite number of such invariants. In particular, this holds for any representation of a compact group.

**Theorem 2.2** (Hilbert) *Let  $\Gamma$  be a compact group and  $V$  a finite-dimensional  $\Gamma$ -module. Then there exist finitely many polynomial invariants  $f_1, \dots, f_{N_{\text{inv}}} : V \rightarrow \mathbb{R}$  such that any polynomial invariant  $r : V \rightarrow \mathbb{R}$  can be expressed as*

$$r(\mathbf{x}) = \tilde{r}(f_1(\mathbf{x}), \dots, f_{N_{\text{inv}}}(\mathbf{x}))$$

with some polynomial  $\tilde{r}$  of  $N_{\text{inv}}$  variables.

See, e.g., Kraft and Procesi [21] for a modern expositions of the invariant theory and Hilbert’s theorem. We refer to the set  $\{f_s\}_{s=1}^{N_{\text{inv}}}$  from this theorem as a *generating set* of polynomial invariants (note that this set is not unique and  $N_{\text{inv}}$  may be different for different generating sets).

Thanks to the density of polynomials in the space of continuous functions, we can easily combine Hilbert’s theorem with the universal approximation theorem to obtain a complete invariant ansatz for invariant maps:

**Proposition 2.3** *Let  $\Gamma$  be a compact group,  $V$  a finite-dimensional  $\Gamma$ -module, and  $f_1, \dots, f_{N_{\text{inv}}} : V \rightarrow \mathbb{R}$  a finite generating set of polynomial invariants on  $V$  (existing by Hilbert’s theorem). Then, any continuous invariant map  $f : V \rightarrow \mathbb{R}$  can be approximated by invariant maps  $\hat{f} : V \rightarrow \mathbb{R}$  of the form*

$$\hat{f}(\mathbf{x}) = \sum_{n=1}^N c_n \sigma \left( \sum_{s=1}^{N_{\text{inv}}} w_{ns} f_s(\mathbf{x}) + h_n \right) \tag{2.5}$$

with some parameter  $N$  and coefficients  $c_n, w_{ns}, h_n$ .

**Proof** It is obvious that the expressions  $\hat{f}$  are  $\Gamma$ -invariant, so we only need to prove the completeness part.

Let us first show that the map  $f$  can be approximated by an invariant polynomial. Let  $K$  be a compact subset in  $V$ , and, like before, consider the symmetrized set  $K_{\text{sym}}$ . By the Stone-Weierstrass theorem, for any  $\epsilon > 0$  there exists a polynomial  $r$  on  $V$  such that  $|r(\mathbf{x}) - f(\mathbf{x})| \leq \epsilon$  for  $\mathbf{x} \in K_{\text{sym}}$ . Consider the symmetrized function  $r_{\text{sym}}(\mathbf{x}) = \int_{\Gamma} r(R_{\gamma}\mathbf{x})d\gamma$ . Then the function  $r_{\text{sym}}$  is invariant and  $|r_{\text{sym}}(\mathbf{x}) - f(\mathbf{x})| \leq \epsilon$  for  $\mathbf{x} \in K$ . On the other hand,  $r_{\text{sym}}$  is a polynomial, since  $r(R_{\gamma}\mathbf{x})$  is a fixed degree polynomial in  $\mathbf{x}$  for any  $\gamma$ .

Using Hilbert’s theorem, we express  $r_{\text{sym}}(\mathbf{x}) = \tilde{r}(f_1(\mathbf{x}), \dots, f_{N_{\text{inv}}}(\mathbf{x}))$  with some polynomial  $\tilde{r}$ .

It remains to approximate the polynomial  $\tilde{r}(z_1, \dots, z_{N_{\text{inv}}})$  by an expression of the form  $\hat{f}(z_1, \dots, z_{N_{\text{inv}}}) = \sum_{n=1}^N c_n \sigma(\sum_{s=1}^{N_{\text{inv}}} w_{ns} z_s + h_n)$  on the compact set  $\{(f_1(\mathbf{x}), \dots, f_{N_{\text{inv}}}(\mathbf{x})) | \mathbf{x} \in K\} \subset \mathbb{R}^{N_{\text{inv}}}$ . By Theorem 2.1, we can do it with any accuracy  $\epsilon$ . Setting finally  $\hat{f}(\mathbf{x}) = \hat{f}(f_1(\mathbf{x}), \dots, f_{N_{\text{inv}}}(\mathbf{x}))$ , we obtain  $\hat{f}$  of the required form such that  $|\hat{f}(\mathbf{x}) - f(\mathbf{x})| \leq 2\epsilon$  for all  $\mathbf{x} \in K$ . □

Note that Proposition 2.3 is a generalization of Theorem 2.1; the latter is a special case obtained if the group is trivial ( $\Gamma = \{e\}$ ) or its representation is trivial ( $R_\gamma \mathbf{x} \equiv \mathbf{x}$ ), and in this case we can just take  $N_{\text{inv}} = d$  and  $f_s(\mathbf{x}) = x_s$ .

In terms of neural network architectures, formula (2.5) can be viewed as a shallow neural network with an extra polynomial layer that precedes the conventional linear combination and nonlinear activation layers.

We extend now the obtained result to equivariant maps. Given two  $\Gamma$ -modules  $V$  and  $U$ , we say that a map  $f : V \rightarrow U$  is *polynomial* if  $l \circ f$  is a polynomial for any linear functional  $l : U \rightarrow \mathbb{R}$ . We rely on the extension of Hilbert’s theorem to polynomial equivariants:

**Lemma 2.1** *Let  $\Gamma$  be a compact group and  $V$  and  $U$  two finite-dimensional  $\Gamma$ -modules. Then there exist finitely many polynomial invariants  $f_1, \dots, f_{N_{\text{inv}}} : V \rightarrow \mathbb{R}$  and polynomial equivariants  $g_1, \dots, g_{N_{\text{eq}}} : V \rightarrow U$  such that any polynomial equivariant  $r_{\text{sym}} : V \rightarrow U$  can be represented in the form  $r_{\text{sym}}(\mathbf{x}) = \sum_{m=1}^{N_{\text{eq}}} g_m(\mathbf{x}) \tilde{r}_m(f_1(\mathbf{x}), \dots, f_{N_{\text{inv}}}(\mathbf{x}))$  with some polynomials  $\tilde{r}_m$ .*

**Proof** We give a sketch of the proof, see e.g. Section 4 of Worfolk [47] for details. A polynomial equivariant  $r_{\text{sym}} : V \rightarrow U$  can be viewed as an invariant element of the space  $\mathbb{R}[V] \otimes U$  with the naturally induced action of  $\Gamma$ , where  $\mathbb{R}[V]$  denotes the space of polynomials on  $V$ . The space  $\mathbb{R}[V] \otimes U$  is in turn a subspace of the algebra  $\mathbb{R}[V \oplus U^*]$ , where  $U^*$  denotes the dual of  $U$ . By Hilbert’s theorem, all invariant elements in  $\mathbb{R}[V \oplus U^*]$  can be generated as polynomials of finitely many invariant elements of this algebra. The algebra  $\mathbb{R}[V \oplus U^*]$  is graded by the degree of the  $U^*$  component, and the corresponding decomposition of  $\mathbb{R}[V \oplus U^*]$  into the direct sum of  $U^*$ -homogeneous spaces indexed by the  $U^*$ -degree  $d_{U^*} = 0, 1, \dots$ , is preserved by the group action. The finitely many polynomials generating all invariant polynomials in  $\mathbb{R}[V \oplus U^*]$  can also be assumed to be  $U^*$ -homogeneous. Let  $\{f_s\}_{s=1}^{N_{\text{inv}}}$  be those of these generating polynomials with  $d_{U^*} = 0$  and  $\{g_s\}_{s=1}^{N_{\text{eq}}}$  be those with  $d_{U^*} = 1$ . Then, a polynomial in the generating invariants is  $U^*$ -homogeneous with  $d_{U^*} = 1$  if and only if it is a linear combination of monomials  $g_s f_1^{n_1} f_2^{n_2} \dots f_{N_{\text{inv}}}^{n_{N_{\text{inv}}}}$ . This yields the representation stated in the lemma.  $\square$

We will refer to the set  $\{g_s\}_{s=1}^{N_{\text{eq}}}$  as a generating set of polynomial equivariants.

The equivariant analog of Proposition 2.3 now reads:

**Proposition 2.4** *Let  $\Gamma$  be a compact group,  $V$  and  $U$  be two finite-dimensional  $\Gamma$ -modules. Let  $f_1, \dots, f_{N_{\text{inv}}} : V \rightarrow \mathbb{R}$  be a finite generating set of polynomial invariants and  $g_1, \dots, g_{N_{\text{eq}}} : V \rightarrow U$  be a finite generating set of polynomial equivariants (existing by Lemma 2.1). Then, any continuous equivariant map  $f : V \rightarrow U$  can be approximated by equivariant maps  $\hat{f} : V \rightarrow U$  of the form*

$$\hat{f}(\mathbf{x}) = \sum_{n=1}^N \sum_{m=1}^{N_{\text{eq}}} c_{mn} g_m(\mathbf{x}) \sigma \left( \sum_{s=1}^{N_{\text{inv}}} w_{mns} f_s(\mathbf{x}) + h_{mn} \right)$$

with some parameter  $N$  and coefficients  $c_{mn}, w_{mns}, h_{mn}$ .

**Proof** The proof is similar to the proof of Proposition 2.3, with the difference that the polynomial map  $r$  is now vector-valued, its symmetrization is defined by  $r_{\text{sym}}(\mathbf{x}) = \int_{\Gamma} R_{\gamma}^{-1} r(R_{\gamma} \mathbf{x}) d\gamma$ , and Lemma 2.1 is used in place of Hilbert’s theorem.  $\square$

We remark that, in turn, Proposition 2.4 generalizes Proposition 2.3; the latter is a special case obtained when  $U = \mathbb{R}$ , and in this case we just take  $N_{\text{eq}} = 1$  and  $g_1 = 1$ .

### 2.3 Polarization and Multiplicity Reduction

The main point of Propositions 2.3 and 2.4 is that the representations described there use *finite* generating sets of invariants and equivariants  $\{f_s\}_{s=1}^{N_{\text{inv}}}$ ,  $\{g_m\}_{m=1}^{N_{\text{eq}}}$  independent of the function  $f$  being approximated. However, the obvious drawback of these results is their non-constructive nature with regard to the functions  $f_s, g_m$ . In general, finding generating sets is not easy. Moreover, the sizes  $N_{\text{inv}}, N_{\text{eq}}$  of these sets in general grow rapidly with the dimensions of the spaces  $V, U$ .

This issue can be somewhat ameliorated using polarization and Weyl’s theorem. Suppose that a  $\Gamma$ -module  $V$  admits a decomposition into a direct sum of invariant submodules:

$$V = \bigoplus_{\alpha} V_{\alpha}^{m_{\alpha}}. \tag{2.6}$$

Here,  $V_{\alpha}^{m_{\alpha}}$  is a direct sum of  $m_{\alpha}$  submodules isomorphic to  $V_{\alpha}$ :

$$V_{\alpha}^{m_{\alpha}} = V_{\alpha} \otimes \mathbb{R}^{m_{\alpha}} = \underbrace{V_{\alpha} \oplus \dots \oplus V_{\alpha}}_{m_{\alpha}}. \tag{2.7}$$

Any finite-dimensional representation of a compact group is completely reducible and has a decomposition of the form (2.6) with non-isomorphic irreducible submodules  $V_{\alpha}$ . In this case the decomposition (2.6) is referred to as the *isotypic decomposition*, and the subspaces  $V_{\alpha}^{m_{\alpha}}$  are known as *isotypic components*. Such isotypic components and their multiplicities  $m_{\alpha}$  are uniquely determined (though individually, the  $m_{\alpha}$  spaces  $V_{\alpha}$  appearing in the direct sum (2.7) are not uniquely determined in general, as subspaces in  $V$ ).

For finite groups the number of non-isomorphic irreducibles  $\alpha$  is finite. In this case, if the module  $V$  is high-dimensional, then this necessarily means that (some of) the multiplicities  $m_{\alpha}$  are large. This is not so, in general, for infinite groups, since infinite compact groups have countably many non-isomorphic irreducible representations. Nevertheless, it is in any case useful to simplify the structure of invariants for high-multiplicity modules, which is what polarization and Weyl’s theorem do.

Below, we slightly abuse the terminology and speak of isotypic components and decompositions in the broader sense, assuming decompositions (2.6), (2.7) but not requiring the submodules  $V_{\alpha}$  to be irreducible or mutually non-isomorphic.

The idea of polarization is to generate polynomial invariants of a representation with large multiplicities from invariants of a representation with small multiplicities. Namely, note that in each isotypic component  $V_{\alpha}^{m_{\alpha}}$  written as  $V_{\alpha} \otimes \mathbb{R}^{m_{\alpha}}$  the group

essentially acts only on the first factor,  $V_\alpha$ . So, given two isotypic  $\Gamma$ -modules of the same type,  $V_\alpha^{m_\alpha} = V_\alpha \otimes \mathbb{R}^{m_\alpha}$  and  $V_\alpha^{m'_\alpha} = V_\alpha \otimes \mathbb{R}^{m'_\alpha}$ , the group action commutes with any linear map  $\mathbb{1}_{V_\alpha} \otimes A : V_\alpha^{m_\alpha} \rightarrow V_\alpha^{m'_\alpha}$ , where  $A$  acts on the second factor,  $A : \mathbb{R}^{m_\alpha} \rightarrow \mathbb{R}^{m'_\alpha}$ . Consequently, given two modules  $V = \bigoplus_\alpha V_\alpha^{m_\alpha}$ ,  $V' = \bigoplus_\alpha V_\alpha^{m'_\alpha}$  and a linear map  $A_\alpha : V_\alpha^{m_\alpha} \rightarrow V_\alpha^{m'_\alpha}$  for each  $\alpha$ , the linear operator  $\mathbf{A} : V \rightarrow V'$  defined by

$$\mathbf{A} = \bigoplus_\alpha \mathbb{1}_{V_\alpha} \otimes A_\alpha \tag{2.8}$$

will commute with the group action. In particular, if  $f$  is a polynomial invariant on  $V'$ , then  $f \circ \mathbf{A}$  will be a polynomial invariant on  $V$ .

The fundamental theorem of Weyl states that it suffices to take  $m'_\alpha = \dim V_\alpha$  to generate in this way a complete set of invariants for  $V$ . We will state this theorem in the following form suitable for our purposes.

**Theorem 2.3** (Weyl [46], sections II.4-5) *Let  $F$  be the set of polynomial invariants for a  $\Gamma$ -module  $V' = \bigoplus_\alpha V_\alpha^{\dim V_\alpha}$ . Suppose that a  $\Gamma$ -module  $V$  admits a decomposition  $V = \bigoplus_\alpha V_\alpha^{m_\alpha}$  with the same  $V_\alpha$ , but arbitrary multiplicities  $m_\alpha$ . Then the polynomials  $\{f \circ \mathbf{A}\}_{f \in F}$  linearly span the space of polynomial invariants on  $V$ , i.e. any polynomial invariant  $f$  on  $V$  can be expressed as  $f(\mathbf{x}) = \sum_{t=1}^T f_t(\mathbf{A}_t \mathbf{x})$  with some polynomial invariants  $f_t$  on  $V'$ .*

**Proof** A detailed exposition of polarization and a proof of Weyl’s theorem based on the Capelli–Deruyts expansion can be found in Weyl’s book or in Sections 7–9 of Kraft and Procesi [21]. We sketch the main idea of the proof.

Consider first the case where  $V$  has only one isotypic component:  $V = V_\alpha^{m_\alpha}$ . We may assume without loss of generality that  $m_\alpha > \dim V_\alpha$  (otherwise the statement is trivial). It is also convenient to identify the space  $V' = V_\alpha^{\dim V_\alpha}$  with the subspace of  $V$  spanned by the first  $\dim V_\alpha$  components  $V_\alpha$ . It suffices to establish the claimed expansion for polynomials  $f$  multihomogeneous with respect to the decomposition  $V = V_\alpha \oplus \dots \oplus V_\alpha$ , i.e. homogeneous with respect to each of the  $m_\alpha$  components. For any such polynomial, the Capelli–Deruyts expansion represents  $f$  as a finite sum  $f = \sum_n C_n B_n f$ . Here  $C_n, B_n$  are linear operators on the space of polynomials on  $V$ , and they belong to the algebra generated by polarization operators on  $V$ . Moreover, for each  $n$ , the polynomial  $\tilde{f}_n = B_n f$  depends only on variables from the first  $\dim V_\alpha$  components of  $V = V_\alpha^{m_\alpha}$ , i.e.  $\tilde{f}_n$  is a polynomial on  $V'$ . This polynomial is invariant, since polarization operators commute with the group action. Since  $C_n$  belongs to the algebra generated by polarization operators, we can then argue (see Proposition 7.4 in Kraft and Procesi [21]) that  $C_n B_n f$  can be represented as a finite sum  $C_n B_n f(\mathbf{x}) = \sum_k \tilde{f}_n((\mathbb{1}_{V_\alpha} \otimes A_{kn})\mathbf{x})$  with some  $m_\alpha \times \dim V_\alpha$  matrices  $A_{kn}$ . This implies the claim of the theorem in the case of a single isotypic component.

Generalization to several isotypic components is obtained by iteratively applying the Capelli–Deruyts expansion to each component. □

Now we can give a more constructive version of Proposition 2.3:

**Proposition 2.5** *Let  $(f_s)_{s=1}^{N_{\text{inv}}}$  be a generating set of polynomial invariants for a  $\Gamma$ -module  $V' = \bigoplus_{\alpha} V_{\alpha}^{\dim V_{\alpha}}$ . Suppose that a  $\Gamma$ -module  $V$  admits a decomposition  $V = \bigoplus_{\alpha} V_{\alpha}^{m_{\alpha}}$  with the same  $V_{\alpha}$ , but arbitrary multiplicities  $m_{\alpha}$ . Then any continuous invariant map  $f : V \rightarrow \mathbb{R}$  can be approximated by invariant maps  $\widehat{f} : V \rightarrow \mathbb{R}$  of the form*

$$\widehat{f}(\mathbf{x}) = \sum_{t=1}^T c_t \sigma \left( \sum_{s=1}^{N_{\text{inv}}} w_{st} f_s(\mathbf{A}_t \mathbf{x}) + h_t \right) \tag{2.9}$$

with some parameter  $T$  and coefficients  $c_t, w_{st}, h_t, \mathbf{A}_t$ , where each  $\mathbf{A}_t$  is formed by an arbitrary collection of  $(m_{\alpha} \times \dim V_{\alpha})$ -matrices  $A_{\alpha}$  as in (2.8).

**Proof** We follow the proof of Proposition 2.3 and approximate the function  $f$  by an invariant polynomial  $r_{\text{sym}}$  on a compact set  $K_{\text{sym}} \subset V$ . Then, using Theorem 2.3, we represent

$$r_{\text{sym}}(\mathbf{x}) = \sum_{t=1}^T r_t(\mathbf{A}_t \mathbf{x}) \tag{2.10}$$

with some invariant polynomials  $r_t$  on  $V'$ . Then, by Proposition 2.3, for each  $t$  we can approximate  $r_t(\mathbf{y})$  on  $\mathbf{A}_t K_{\text{sym}}$  by an expression

$$\sum_{n=1}^N \widetilde{c}_{nt} \sigma \left( \sum_{s=1}^{N_{\text{inv}}} \widetilde{w}_{nst} f_s(\mathbf{y}) + \widetilde{h}_{nt} \right) \tag{2.11}$$

with some  $\widetilde{c}_{nt}, \widetilde{w}_{nst}, \widetilde{h}_{nt}$ . Combining (2.10) with (2.11), it follows that  $f$  can be approximated on  $K_{\text{sym}}$  by

$$\sum_{t=1}^T \sum_{n=1}^N \widetilde{c}_{nt} \sigma \left( \sum_{s=1}^{N_{\text{inv}}} \widetilde{w}_{nst} f_s(\mathbf{A}_t \mathbf{x}) + \widetilde{h}_{nt} \right).$$

The final expression (2.9) is obtained now by removing the superfluous summation over  $n$ . □

Proposition 2.5 is more constructive than Proposition 2.3 in the sense that the approximating ansatz (2.9) only requires us to know an isotypic decomposition  $V = \bigoplus_{\alpha} V_{\alpha}^{m_{\alpha}}$  of the  $\Gamma$ -module under consideration and a generating set  $(f_s)_{s=1}^{N_{\text{inv}}}$  for the reference module  $V' = \bigoplus_{\alpha} V_{\alpha}^{\dim V_{\alpha}}$ . In particular, suppose that the group  $\Gamma$  is finite, so that there are only finitely many non-isomorphic irreducible modules  $V_{\alpha}$ . Then, for any  $\Gamma$ -module  $V$ , the universal approximating ansatz (2.9) includes not more than  $CT \dim V$  scalar weights, with some constant  $C$  depending only on  $\Gamma$  (since  $\dim V = \sum_{\alpha} m_{\alpha} \dim V_{\alpha}$ ).

We remark that in terms of the network architecture, formula (2.9) can be interpreted as the network (2.5) from Proposition 2.3 with an extra linear layer performing multiplication of the input vector by  $\mathbf{A}_t$ .

We establish now an equivariant analog of Proposition 2.5. We start with an equivariant analog of Theorem 2.3.

**Proposition 2.6** *Let  $V' = \bigoplus_{\alpha} V_{\alpha}^{\dim V_{\alpha}}$  and  $G$  be the space of polynomial equivariants  $g : V' \rightarrow U$ . Suppose that a  $\Gamma$ -module  $V$  admits a decomposition  $V = \bigoplus_{\alpha} V_{\alpha}^{m_{\alpha}}$  with the same  $V_{\alpha}$ , but arbitrary multiplicities  $m_{\alpha}$ . Then, the functions  $\{g \circ \mathbf{A}\}_{g \in G}$  linearly span the space of polynomial equivariants  $g : V \rightarrow U$ , i.e. any such equivariant can be expressed as  $g(\mathbf{x}) = \sum_{t=1}^T g_t(\mathbf{A}_t \mathbf{x})$  with some polynomial equivariants  $g_t : V' \rightarrow U$ .*

**Proof** As mentioned in the proof of Lemma 2.1, polynomial equivariants  $g : V \rightarrow U$  can be viewed as invariant elements of the extended polynomial algebra  $\mathbb{R}[V \oplus U^*]$ . The proof of the theorem is then completely analogous to the proof of Theorem 2.3 and consists in applying the Capelli–Deruyts expansion to each isotypic component of the submodule  $V$  in  $V \oplus U^*$ .  $\square$

The equivariant analog of Proposition 2.5 now reads:

**Proposition 2.7** *Let  $(f_s)_{s=1}^{N_{\text{inv}}}$  be a generating set of polynomial invariants for a  $\Gamma$ -module  $V' = \bigoplus_{\alpha} V_{\alpha}^{\dim V_{\alpha}}$ , and  $(g_s)_{s=1}^{N_{\text{eq}}}$  be a generating sets of polynomial equivariants mapping  $V'$  to a  $\Gamma$ -module  $U$ . Let  $V = \bigoplus_{\alpha} V_{\alpha}^{m_{\alpha}}$  be a  $\Gamma$ -module with the same  $V_{\alpha}$ . Then any continuous equivariant map  $f : V \rightarrow U$  can be approximated by equivariant maps  $\hat{f} : V \rightarrow U$  of the form*

$$\hat{f}(\mathbf{x}) = \sum_{t=1}^T \sum_{m=1}^{N_{\text{eq}}} c_{mt} g_m(\mathbf{A}_t \mathbf{x}) \sigma \left( \sum_{s=1}^{N_{\text{inv}}} w_{mst} f_s(\mathbf{A}_t \mathbf{x}) + h_{mt} \right) \tag{2.12}$$

with some coefficients  $c_{mt}, w_{mst}, h_{mt}, \mathbf{A}_t$ , where each  $\mathbf{A}_t$  is given by a collection of  $(m_{\alpha} \times \dim V_{\alpha})$ -matrices  $A_{\alpha}$  as in (2.8).

**Proof** As in the proof of Theorem 2.4, we approximate the function  $f$  by a polynomial equivariant  $r_{\text{sym}}$  on a compact  $K_{\text{sym}} \subset V$ . Then, using Theorem 2.6, we represent

$$r_{\text{sym}}(\mathbf{x}) = \sum_{t=1}^T r_t(\mathbf{A}_t \mathbf{x}) \tag{2.13}$$

with some polynomial equivariants  $r_t : V' \rightarrow U$ . Then, by Proposition 2.4, for each  $t$  we can approximate  $r_t(\mathbf{x}')$  on  $\mathbf{A}_t K_{\text{sym}}$  by expressions

$$\sum_{n=1}^N \sum_{m=1}^{N_{\text{eq}}} \tilde{c}_{mnt} g(\mathbf{x}') \sigma \left( \sum_{s=1}^{N_{\text{inv}}} \tilde{w}_{mnst} f_s(\mathbf{x}') + \tilde{h}_{mnt} \right). \tag{2.14}$$

Using (2.13) and (2.14),  $f$  can be approximated on  $K_{\text{sym}}$  by expressions

$$\sum_{t=1}^T \sum_{n=1}^N \sum_{m=1}^{N_{\text{eq}}} \tilde{c}_{mnt} g(\mathbf{A}_t \mathbf{x}) \sigma \left( \sum_{s=1}^{N_{\text{inv}}} \tilde{w}_{mnst} f_s(\mathbf{A}_t \mathbf{x}) + \tilde{h}_{mnt} \right).$$

We obtain the final form (2.12) by removing the superfluous summation over  $n$ .  $\square$

We remark that Proposition 2.7 improves the earlier Proposition 2.4 in the equivariant setting in the same sense in which Proposition 2.5 improves Proposition 2.3 in the invariant setting: construction of a universal approximator in the case of arbitrary isotypic multiplicities is reduced to the construction with particular multiplicities by adding an extra equivariant linear layer to the network.

### 2.4 The Symmetric Group $S_N$

Even with the simplification resulting from polarization, the general results of the previous section are not immediately useful, since one still needs to find the isotypic decomposition of the analyzed  $\Gamma$ -modules and to find the relevant generating invariants and equivariants. In this section we describe one particular case where the approximating expression can be reduced to a fully explicit form.

Namely, consider the natural action of the symmetric group  $S_N$  on  $\mathbb{R}^N$ :

$$R_\gamma \mathbf{e}_n = \mathbf{e}_{\gamma(n)},$$

where  $\mathbf{e}_n \in \mathbb{R}^N$  is a coordinate vector and  $\gamma \in S_N$  is a permutation.

Let  $V = \mathbb{R}^N \otimes \mathbb{R}^M$  and consider  $V$  as a  $S_N$ -module by assuming that the group acts on the first factor, i.e.  $\gamma$  acts on  $\mathbf{x} = \sum_{n=1}^N \mathbf{e}_n \otimes \mathbf{x}_n \in V$  by

$$R_\gamma \sum_{n=1}^N \mathbf{e}_n \otimes \mathbf{x}_n = \sum_{n=1}^N \mathbf{e}_{\gamma(n)} \otimes \mathbf{x}_n.$$

We remark that this module appears, for example, in the following scenario (cf. Zaheer et al. [49]). Suppose that  $f$  is a map defined on the set of sets  $X = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$  of  $N$  vectors from  $\mathbb{R}^M$ . We can identify the set  $X$  with the element  $\sum_{n=1}^N \mathbf{e}_n \otimes \mathbf{x}_n$  of  $V$  and in this way view  $f$  as defined on a subset of  $V$ . However, since the set  $X$  is unordered, it can also be identified with  $\sum_{n=1}^N \mathbf{e}_{\gamma(n)} \otimes \mathbf{x}_n$  for any permutation  $\gamma \in S_N$ . Accordingly, if the map  $f$  is to be extended to the whole  $V$ , then this extension needs to be invariant with respect to the above action of  $S_N$ .

We describe now an explicit complete ansatz for  $S_N$ -invariant approximations of functions on  $V$ . This is made possible by another classical theorem of Weyl and by a simple form of a generating set of permutation invariants on  $\mathbb{R}^N$ . We will denote by



$x_{nm}$  the coordinates of  $\mathbf{x} \in V$  with respect to the canonical basis in  $V$ :

$$\mathbf{x} = \sum_{n=1}^N \sum_{m=1}^M x_{nm} \mathbf{e}_n \otimes \mathbf{e}_m.$$

**Theorem 2.4** *Let  $V = \mathbb{R}^N \otimes \mathbb{R}^M$  and  $f : V \rightarrow \mathbb{R}$  be a  $S_N$ -invariant continuous map. Then  $f$  can be approximated by  $S_N$ -invariant expressions*

$$\widehat{f}(\mathbf{x}) = \sum_{t=1}^{T_1} c_t \sigma \left( \sum_{q=1}^{T_2} w_{qt} \sum_{n=1}^N \sigma \left( b_q \sum_{m=1}^M a_{tm} x_{nm} + e_q \right) + h_t \right), \tag{2.15}$$

with some parameters  $T_1, T_2$  and coefficients  $c_t, w_{qt}, b_q, a_{tm}, e_q, h_t$ .

**Proof** It is clear that expression (2.15) is  $S_N$ -invariant and we only need to prove its completeness. The theorem of Weyl [46, Section II.3] states that a generating set of symmetric polynomials on  $V$  can be obtained by polarizing a generating set of symmetric polynomials  $\{f_p\}_{p=1}^{N_{\text{inv}}}$  defined on a *single* copy of  $\mathbb{R}^N$ . Arguing as in Proposition 2.5, it follows that any  $S_N$ -invariant continuous map  $f : V \rightarrow \mathbb{R}$  can be approximated by expressions

$$\sum_{t=1}^{T_1} \tilde{c}_t \sigma \left( \sum_{p=1}^{N_{\text{inv}}} \tilde{w}_{pt} f_p(\tilde{\mathbf{A}}_t \mathbf{x}) + \tilde{h}_t \right),$$

where  $\tilde{\mathbf{A}}_t \mathbf{x} = \sum_{n=1}^N \sum_{m=1}^M \tilde{a}_{tm} x_{nm} \mathbf{e}_n$ . A well-known generating set of symmetric polynomials on  $\mathbb{R}^N$  is the first  $N$  coordinate power sums:

$$f_p(\mathbf{y}) = \sum_{n=1}^N \tilde{f}_p(y_n), \text{ where } \mathbf{y} = (y_1, \dots, y_N), \tilde{f}_p(y_n) = y_n^p, \quad p = 1, \dots, N.$$

It follows that  $f$  can be approximated by expressions

$$\sum_{t=1}^{T_1} \tilde{c}_t \sigma \left( \sum_{p=1}^N \tilde{w}_{pt} \sum_{n=1}^N \tilde{f}_p \left( \sum_{m=1}^M \tilde{a}_{tm} x_{nm} \right) + \tilde{h}_t \right). \tag{2.16}$$

Using Theorem 2.1, we can approximate  $\tilde{f}_p(y)$  by expressions  $\sum_{q=1}^T \tilde{d}_{pq} \sigma(\tilde{b}_{pq} y + \tilde{h}_{pq})$ . It follows that (2.16) can be approximated by

$$\sum_{t=1}^{T_1} \tilde{c}_t \sigma \left( \sum_{p=1}^N \sum_{q=1}^T \tilde{w}_{pt} \tilde{d}_{pq} \sum_{n=1}^N \sigma \left( \tilde{b}_{pq} \sum_{m=1}^M \tilde{a}_{tm} x_{nm} + h_{pq} \right) + \tilde{h}_t \right).$$

Replacing the double summation over  $p, q$  by a single summation over  $q$ , we arrive at (2.15). □

Note that expression (2.15) resembles the formula of the usual (non-invariant) feedforward network with two hidden layers of sizes  $T_1$  and  $T_2$ :

$$\hat{f}(\mathbf{x}) = \sum_{t=1}^{T_1} c_t \sigma \left( \sum_{q=1}^{T_2} w_{qt} \sigma \left( \sum_{n=1}^N \sum_{m=1}^M a_{qnm} x_{nm} + e_q \right) + h_t \right).$$

Let us also compare ansatz (2.15) with the ansatz obtained by direct symmetrization (see Proposition (2.1)), which in our case has the form

$$\hat{f}(\mathbf{x}) = \sum_{\gamma \in S_N} \sum_{t=1}^T c_t \sigma \left( \sum_{n=1}^N \sum_{m=1}^M w_{\gamma(n),m,t} x_{nm} + h_t \right).$$

From the application perspective, since  $|S_N| = N!$ , at large  $N$  this expression has prohibitively many terms and is therefore impractical without subsampling of  $S_N$ , which would break the exact  $S_N$ -invariance. In contrast, ansatz (2.15) is complete, fully  $S_N$ -invariant and involves only  $O(T_1 N(M + T_2))$  arithmetic operations and evaluations of  $\sigma$ .

We remark that another complete permutation-invariant network architecture, using the max function, was given in Qi et al. [34, Theorem 1].

### 3 Translations and Deep Convolutional Networks

Convolutional neural networks (convnets, [22]) play a key role in many modern applications of deep learning. Such networks operate on input data having grid-like structure

(usually, spatial or temporal) and consist of multiple stacked convolutional layers transforming initial object description into increasingly complex features necessary to recognize complex patterns in the data. The shape of earlier layers in the network mimics the shape of input data, but later layers gradually become “thinner” geometrically while acquiring “thicker” *feature dimensions*. We refer the reader to deep learning literature for details on these networks, e.g. see Chapter 9 in Goodfellow et al. [12] for an introduction.

There are several important concepts associated with convolutional networks, in particular *weight sharing* (which ensures approximate *translation equivariance* of the layers with respect to grid shifts); *locality* of the layer operation; and *pooling*. Locality means that the layer output at a certain geometric point of the domain depends only on a small neighborhood of this point. Pooling is a grid subsampling that helps reshape the data flow by removing excessive spatial detailization. Practical usefulness of convnets stems from the interplay between these various elements of convnet design.

From the perspective of the main topic of the present work—group invariant/equivariant networks—we are mostly interested in invariance/equivariance of convnets with respect to Lie groups such as the group of translations or the group of rigid motions (to be considered in Sect. 4), and we would like to establish relevant universal approximation theorems. However, we first point out some serious difficulties that one faces when trying to formulate and prove such results.

**Lack of symmetry in finite computational models** Practically used convnets are finite models; in particular they operate on discretized and bounded domains that do not possess the full symmetry of the spaces  $\mathbb{R}^d$ . While the translational symmetry is partially preserved by discretization to a regular grid, and the group  $\mathbb{R}^d$  can be in a sense approximated by the groups  $(\lambda\mathbb{Z})^d$  or  $(\lambda\mathbb{Z}_n)^d$ , one cannot reconstruct, for example, the rotational symmetry in a similar way. If a group  $\Gamma$  is compact, then, as discussed in Sect. 2, we can still obtain finite and fully  $\Gamma$ -invariant/equivariant computational models by considering finite-dimensional representations of  $\Gamma$ , but this is not the case with noncompact groups such as  $\mathbb{R}^d$ . Therefore, in the case of the group  $\mathbb{R}^d$  (and the group of rigid planar motions considered later in Sect. 4), we will need to prove the desired results on invariance/equivariance and completeness of convnets only in the limit of infinitely large domain and infinitesimal grid spacing.

**Erosion of translation equivariance by pooling** Pooling reduces the translational symmetry of the convnet model. For example, if a few first layers of the network define a map equivariant with respect to the group  $(\lambda\mathbb{Z})^2$  with some spacing  $\lambda$ , then after pooling with stride  $m$  the result will only be equivariant with respect to the subgroup  $(m\lambda\mathbb{Z})^2$ . (We remark in this regard that in practical applications, weight sharing and accordingly translation equivariance are usually only important for earlier layers of convolutional networks.) Therefore, we will consider separately the cases of convnets without or with pooling; the  $\mathbb{R}^d$ -equivariance will only apply in the former case.

In view of the above difficulties, in this section we will give several versions of the universal approximation theorem for convnets, with different treatments of these issues.

In Sect. 3.1 we prove a universal approximation theorem for a single non-local convolutional layer on a finite discrete grid with periodic boundary conditions

(Proposition 3.1). This basic result is a straightforward consequence of the general Proposition 2.2 when applied to finite abelian groups.

In Sect. 3.2 we prove the main result of Sect. 3, Theorem 3.1. This theorem extends Proposition 3.1 in several important ways. First, we will consider continuum signals, i.e. assume that the approximated map is defined on functions on  $\mathbb{R}^n$  rather than on functions on a discrete grid. This extension will later allow us to rigorously formulate a universal approximation theorem for rotations and euclidean motions in Sect. 4. Second, we will consider stacked convolutional layers and assume each layer to act locally (as in convnets actually used in applications). However, the setting of Theorem 3.2 will not involve pooling, since, as remarked above, pooling destroys the translation equivariance of the model.

In Sect. 3.3 we prove Theorem 3.2, relevant for convnets most commonly used in practice. Compared to the setting of Sect. 3.2, this computational model will be spatially bounded, will include pooling, and will not assume translation invariance of the approximated map.

### 3.1 Finite Abelian Groups and Single Convolutional Layers

We consider a group

$$\Gamma = \mathbb{Z}_{n_1} \times \cdots \times \mathbb{Z}_{n_v}, \quad (3.1)$$

where  $\mathbb{Z}_n = \mathbb{Z}/(n\mathbb{Z})$  is the cyclic group of order  $n$ . Note that the group  $\Gamma$  is abelian and conversely, by the fundamental theorem of finite abelian groups, any such group can be represented in the form (3.1).

We consider the “input” module  $V = \mathbb{R}^\Gamma \otimes \mathbb{R}^{d_V}$  and the “output” module  $U = \mathbb{R}^\Gamma \otimes \mathbb{R}^{d_U}$ , with some finite dimensions  $d_V, d_U$  and with the natural representation of  $\Gamma$ :

$$R_\gamma(\mathbf{e}_\theta \otimes \mathbf{v}) = \mathbf{e}_{\theta+\gamma} \otimes \mathbf{v}, \quad \gamma, \theta \in \Gamma, \quad \mathbf{v} \in \mathbb{R}^{d_V} \text{ or } \mathbb{R}^{d_U}.$$

We will denote elements of  $V, U$  by boldface characters  $\Phi$  and interpret them as  $d_V$ - or  $d_U$ -component *signals* defined on the set  $\Gamma$ . For example, in the context of 2D image processing we have  $v = 2$  and the group  $\Gamma = \mathbb{Z}_{n_1} \times \mathbb{Z}_{n_2}$  corresponds to a discretized rectangular image with periodic boundary conditions, where  $n_1, n_2$  are the geometric sizes of the image while  $d_V$  and  $d_U$  are the numbers of input and output features, respectively (in particular, if the input is a usual RGB image, then  $d_V = 3$ ).

Denote by  $\Phi_{\theta k}$  the coefficients in the expansion of a vector  $\Phi$  from  $V$  or  $U$  over the standard product bases in these spaces:

$$\Phi = \sum_{\theta \in \Gamma} \sum_{k=1}^{d_V \text{ or } d_U} \Phi_{\theta k} \mathbf{e}_\theta \otimes \mathbf{e}_k. \quad (3.2)$$

We describe now a complete equivariant ansatz for approximating  $\Gamma$ -equivariant maps  $f : V \rightarrow U$ . Thanks to decomposition (2.4), we may assume without loss

that  $d_U = 1$ . By (3.2), any map  $f : V \rightarrow U$  is then specified by the coefficients  $f(\Phi)_\theta (\equiv f(\Phi)_{\theta,1}) \in \mathbb{R}$  as  $\Phi$  runs over  $V$  and  $\theta$  runs over  $\Gamma$ .

**Proposition 3.1** *Any continuous  $\Gamma$ -equivariant map  $f : V \rightarrow U$  can be approximated by  $\Gamma$ -equivariant maps  $\hat{f} : V \rightarrow U$  of the form*

$$\hat{f}(\Phi)_\gamma = \sum_{n=1}^N c_n \sigma \left( \sum_{\theta \in \Gamma} \sum_{k=1}^{d_V} w_{n\theta k} \Phi_{\gamma+\theta,k} + h_n \right), \tag{3.3}$$

where  $\Phi = \sum_{\gamma \in \Gamma} \sum_{k=1}^{d_V} \Phi_{\gamma k} \mathbf{e}_\gamma \otimes \mathbf{e}_k$ ,  $N$  is a parameter, and  $c_n, w_{n\theta k}, h_n$  are some coefficients.

**Proof** We apply Proposition 2.2 with  $l_n(\Phi) = \sum_{\theta' \in \Gamma} \sum_{k=1}^{d_V} w'_{n\theta'k} \Phi_{\theta'k}$  and  $\mathbf{y}_n = \sum_{\varkappa \in \Gamma} y_{n\varkappa} \mathbf{e}_\varkappa$ , and obtain the ansatz

$$\hat{f}(\Phi) = \sum_{\gamma' \in \Gamma} \sum_{n=1}^N \sum_{\varkappa \in \Gamma} y_{n\varkappa} \sigma \left( \sum_{\theta' \in \Gamma} \sum_{k=1}^{d_V} w'_{n\theta'k} \Phi_{\theta'-\gamma',k} + h_n \right) \mathbf{e}_{\varkappa-\gamma'} = \sum_{\varkappa \in \Gamma} \sum_{n=1}^N y_{n\varkappa} \mathbf{a}_{\varkappa n},$$

where

$$\mathbf{a}_{\varkappa n} = \sum_{\gamma' \in \Gamma} \sigma \left( \sum_{\theta' \in \Gamma} \sum_{k=1}^{d_V} w'_{n\theta'k} \Phi_{\theta'-\gamma',k} + h_n \right) \mathbf{e}_{\varkappa-\gamma'}. \tag{3.4}$$

By linearity of the expression on the r.h.s. of (3.3), it suffices to check that each  $\mathbf{a}_{\varkappa n}$  can be written in the form

$$\sum_{\gamma \in \Gamma} \sigma \left( \sum_{\theta \in \Gamma} \sum_{k=1}^{d_V} w_{n\theta k} \Phi_{\theta+\gamma,k} + h_n \right) \mathbf{e}_\gamma.$$

But this expression results if we make in (3.4) the substitutions  $\gamma = \varkappa - \gamma', \theta = \theta' - \varkappa$  and  $w_{n\theta k} = w'_{n,\theta+\varkappa,k}$ . □

The expression (3.3) resembles the standard convolutional layer without pooling as described, e.g., in Goodfellow et al. [12]. Specifically, this expression can be viewed as a linear combination of  $N$  scalar filters obtained as compositions of linear convolutions with pointwise non-linear activations. An important difference with the standard convolutional layers is that the convolutions in (3.3) are non-local, in the sense that the weights  $w_{n\theta k}$  do not vanish at large  $\theta$ . Clearly, this non-locality is inevitable if approximation is to be performed with just a single convolutional layer.

We remark that it is possible to use Proposition 2.4 to describe an alternative complete  $\Gamma$ -equivariant ansatz based on polynomial invariants and equivariants. However, this approach seems to be less efficient because it is relatively difficult to specify a small explicit set of generating polynomials for abelian groups (see, e.g. Schmid [36] for a number of relevant results). Nevertheless, we will use polynomial invariants of the abelian group  $SO(2)$  in our construction of “charge-conserving convnet” in Sect. 4.

### 3.2 Continuum Signals and Deep Convnets

In this section we extend Proposition 3.1 in several ways.

First, instead of the group  $\mathbb{Z}_{n_1} \times \dots \times \mathbb{Z}_{n_v}$ , we consider the group  $\Gamma = \mathbb{R}^v$ . Accordingly, we will consider infinite-dimensional  $\mathbb{R}^v$ -modules

$$\begin{aligned} V &= L^2(\mathbb{R}^v) \otimes \mathbb{R}^{d_V} \cong L^2(\mathbb{R}^v, \mathbb{R}^{d_V}), \\ U &= L^2(\mathbb{R}^v) \otimes \mathbb{R}^{d_U} \cong L^2(\mathbb{R}^v, \mathbb{R}^{d_U}) \end{aligned}$$

with some finite  $d_V, d_U$ . Here,  $L^2(\mathbb{R}^v, \mathbb{R}^d)$  is the Hilbert space of maps  $\Phi : \mathbb{R}^v \rightarrow \mathbb{R}^d$  with  $\int_{\mathbb{R}^d} |\Phi(\gamma)|^2 d\gamma < \infty$ , equipped with the standard scalar product  $\langle \Phi, \Psi \rangle = \int_{\mathbb{R}^d} \Phi(\gamma) \cdot \Psi(\gamma) d\gamma$ , where  $\Phi(\gamma) \cdot \Psi(\gamma)$  denotes the scalar product of  $\Phi(\gamma)$  and  $\Psi(\gamma)$  in  $\mathbb{R}^d$ . The group  $\mathbb{R}^v$  is naturally represented on  $V, U$  by

$$R_\gamma \Phi(\theta) = \Phi(\theta - \gamma), \quad \Phi \in V \text{ or } U, \quad \gamma, \theta \in \mathbb{R}^v. \tag{3.5}$$

Throughout this section,  $R_\gamma$  will denote this representation of  $\mathbb{R}^v$  on  $V$  or  $U$ . Compared to the setting of the previous subsection, we interpret the modules  $V, U$  as carrying now “infinitely extended” and “infinitely detailed”  $d_V$ - or  $d_U$ -component signals. We will be interested in approximating arbitrary  $\mathbb{R}^v$ -equivariant continuous maps  $f : V \rightarrow U$ .

The second extension is that we will perform this approximation using stacked convolutional layers with local action. Our approximation will be a finite computational model, and to define it we first need to apply a discretization and a spatial cutoff to vectors from  $V$  and  $U$ .

Let us first describe the discretization. For any *grid spacing*  $\lambda > 0$ , let  $V_\lambda$  be the subspace in  $V$  formed by signals  $\Phi : \mathbb{R}^v \rightarrow \mathbb{R}^{d_V}$  constant on all cubes

$$Q_{\mathbf{k}}^{(\lambda)} = \prod_{s=1}^v \left[ (k_s - \frac{1}{2})\lambda, (k_s + \frac{1}{2})\lambda \right],$$

where  $\mathbf{k} = (k_1, \dots, k_v) \in \mathbb{Z}^v$ . Let  $P_\lambda$  be the orthogonal projector onto  $V_\lambda$  in  $V$ :

$$P_\lambda \Phi(\gamma) = \frac{1}{\lambda^v} \int_{Q_{\mathbf{k}}^{(\lambda)}} \Phi(\theta) d\theta, \quad \text{where } Q_{\mathbf{k}}^{(\lambda)} \ni \gamma. \tag{3.6}$$

A function  $\Phi \in V_\lambda$  can naturally be viewed as a function on the lattice  $(\lambda\mathbb{Z})^v$ , so that we can also view  $V_\lambda$  as a Hilbert space

$$V_\lambda \cong L^2((\lambda\mathbb{Z})^v, \mathbb{R}^{d_V}),$$

with the scalar product  $\langle \Phi, \Psi \rangle = \lambda^v \sum_{\gamma \in (\lambda\mathbb{Z})^v} \Phi(\gamma) \cdot \Psi(\gamma)$ . We define the subspaces  $U_\lambda \subset U$  similarly to the subspaces  $V_\lambda \subset V$ .

Next, we define the spatial cutoff. For an integer  $L \geq 0$  we denote by  $Z_L$  the size- $2L$  cubic subset of the grid  $\mathbb{Z}^v$ :

$$Z_L = \{\mathbf{k} \in \mathbb{Z}^v \mid \|\mathbf{k}\|_\infty \leq L\}, \tag{3.7}$$

where  $\mathbf{k} = (k_1, \dots, k_v) \in \mathbb{Z}^v$  and  $\|\mathbf{k}\|_\infty = \max_{n=1, \dots, v} |k_n|$ . Let  $\lfloor \cdot \rfloor$  denote the standard floor function. For any  $\Lambda \geq 0$  (referred to as the *spatial range* or *cutoff*) we define the subspace  $V_{\lambda, \Lambda} \subset V_\lambda$  by

$$\begin{aligned} V_{\lambda, \Lambda} &= \{\Phi : (\lambda\mathbb{Z})^v \rightarrow \mathbb{R}^{d_v} \mid \Phi(\lambda\mathbf{k}) = 0 \text{ if } \mathbf{k} \notin Z_{\lfloor \Lambda/\lambda \rfloor}\} \\ &\cong \{\Phi : \lambda Z_{\lfloor \Lambda/\lambda \rfloor} \rightarrow \mathbb{R}^{d_v}\} \\ &\cong L^2(\lambda Z_{\lfloor \Lambda/\lambda \rfloor}, \mathbb{R}^{d_v}). \end{aligned} \tag{3.8}$$

Clearly,  $\dim V_{\lambda, \Lambda} = (2\lfloor \Lambda/\lambda \rfloor + 1)^v d_v$ . The subspaces  $U_{\lambda, \Lambda} \subset U_\lambda$  are defined in a similar fashion. We will denote by  $P_{\lambda, \Lambda}$  the linear operators orthogonally projecting  $V$  to  $V_{\lambda, \Lambda}$  or  $U$  to  $U_{\lambda, \Lambda}$ .

In the following, we will assume that the convolutional layers have a finite *receptive field*  $Z_{L_{\text{rf}}}$ —a set of the form (3.7) with some fixed  $L_{\text{rf}} > 0$ .

We can now describe our model of stacked convnets that will be used to approximate maps  $f : V \rightarrow U$  (see Fig. 1). Namely, our approximation will be a composition of the form

$$\widehat{f} : V \xrightarrow{P_{\lambda, \Lambda + (T-1)\lambda L_{\text{rf}}}} V_{\lambda, \Lambda + (T-1)\lambda L_{\text{rf}}} (\equiv W_1) \xrightarrow{\widehat{f}_1} W_2 \xrightarrow{\widehat{f}_2} \dots \xrightarrow{\widehat{f}_T} W_{T+1} (\equiv U_{\lambda, \Lambda}). \tag{3.9}$$

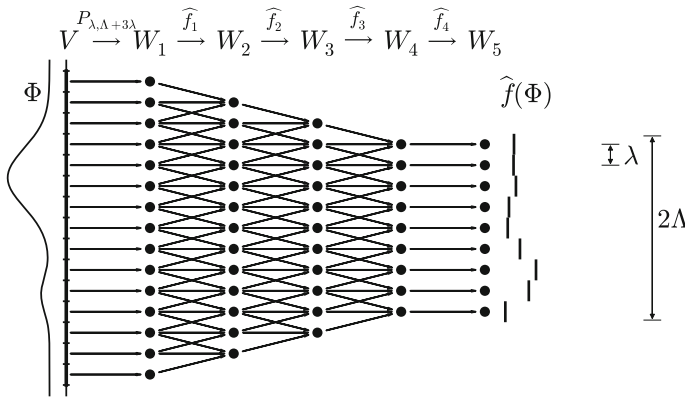
Here, the first step  $P_{\lambda, \Lambda + (T-1)\lambda L_{\text{rf}}}$  is an orthogonal finite-dimensional projection implementing the initial discretization and spatial cutoff of the signal. The maps  $\widehat{f}_t$  are convolutional layers connecting intermediate spaces

$$W_t = \begin{cases} \{\Phi : \lambda Z_{\lfloor \Lambda/\lambda \rfloor + (T-t)L_{\text{rf}}} \rightarrow \mathbb{R}^{d_t}\}, & t \leq T \\ \{\Phi : \lambda Z_{\lfloor \Lambda/\lambda \rfloor} \rightarrow \mathbb{R}^{d_t}\}, & t = T + 1 \end{cases} \tag{3.10}$$

with some *feature dimensions*  $d_t$  such that  $d_1 = d_v$  and  $d_{T+1} = d_u$ . The first intermediate space  $W_1$  is identified with the space  $V_{\lambda, \Lambda + (T-1)\lambda L_{\text{rf}}}$  (the image of the projector  $P_{\lambda, \Lambda + (T-1)\lambda L_{\text{rf}}}$  applied to  $V$ ), while the end space  $W_{T+1}$  is identified with  $U_{\lambda, \Lambda}$  (the respective discretization and cutoff of  $U$ ).

The convolutional layers are defined as follows. Let  $(\Phi_{\gamma n})_{\gamma \in Z_{\lfloor \Lambda/\lambda \rfloor + (T-t)L_{\text{rf}}}, n=1, \dots, d_t}$  be the coefficients in the expansion of  $\Phi \in W_t$  over the standard basis in  $W_t$ , as in (3.2). Then, for  $t < T$  we define  $\widehat{f}_t$  using the conventional “linear convolution followed by nonlinear activation” formula,

$$\widehat{f}_t(\Phi)_{\gamma n} = \sigma \left( \sum_{\theta \in Z_{L_{\text{rf}}}} \sum_{k=1}^{d_t} w_{n\theta k}^{(t)} \Phi_{\gamma + \theta, k} + h_n^{(t)} \right), \quad \gamma \in Z_{\lfloor \Lambda/\lambda \rfloor + (T-t)L_{\text{rf}}}, n = 1, \dots, d_{t+1},$$



**Fig. 1** A one-dimensional ( $v = 1$ ) basic convnet with the receptive field parameter  $L_{\text{rf}} = 1$ . The dots show feature spaces  $\mathbb{R}^{d_t}$  associated with particular points of the grid  $\lambda\mathbb{Z}$

$$(3.11)$$

while in the last layer ( $t = T$ ) we drop nonlinearities and only form a linear combination of values at the same point of the grid:

$$\widehat{f}_T(\Phi)_{\gamma n} = \sum_{k=1}^{d_T} w_{nk}^{(T)} \Phi_{\gamma k} + h_n^{(T)}, \quad \gamma \in Z_{\lfloor \Lambda/\lambda \rfloor}, n = 1, \dots, d_U. \quad (3.12)$$

Note that the grid size  $\lfloor \Lambda/\lambda \rfloor + (T - t)L_{\text{rf}}$  associated with the space  $W_t$  is consistent with the rule (3.11) which evaluates the new signal  $\widehat{f}(\Phi)$  at each node of the grid as a function of the signal  $\Phi$  in the  $L_{\text{rf}}$ -neighborhood of that node (so that the domain  $\lambda Z_{\lfloor \Lambda/\lambda \rfloor + (T-t)L_{\text{rf}}}$  “shrinks” slightly as  $t$  grows).

Note that we can interpret the map  $\widehat{f}$  as a map between  $V$  and  $U$ , since  $U_{\lambda, \Lambda} \subset U$ .

**Definition 3.1** A **basic convnet** is a map  $\widehat{f} : V \rightarrow U$  defined by (3.9), (3.11), (3.12), and characterized by parameters  $\lambda, \Lambda, L_{\text{rf}}, T, d_1, \dots, d_{T+1}$  and coefficients  $w_{n\theta k}^{(t)}$  and  $h_n^{(t)}$ .

Note that, defined in this way, a basic convnet is a finite computational model in the following sense: while being a map between infinite-dimensional spaces  $V$  and  $U$ , all the steps in  $\widehat{f}$  except the initial discretization and cutoff involve only finitely many arithmetic operations and evaluations of the activation function.

We aim to prove an analog of Theorem 2.1, stating that any continuous  $\mathbb{R}^v$ -equivariant map  $f : V \rightarrow U$  can be approximated by basic convnets in the topology of uniform convergence on compact sets. However, there are some important caveats due to the fact that the space  $V$  is now infinite-dimensional.

First, in contrast to the case of finite-dimensional spaces, balls in  $L^2(\mathbb{R}^v, \mathbb{R}^{d_v})$  are not compact. The well-known general criterion states that in a complete metric space, and in particular in  $V = L^2(\mathbb{R}^v, \mathbb{R}^{d_v})$ , a set is compact iff it is closed and *totally bounded*, i.e. for any  $\epsilon > 0$  can be covered by finitely many  $\epsilon$ -balls.



The second point (related to the first) is that a finite-dimensional space is *hemicompact*, i.e., there is a sequence of compact sets such that any other compact set is contained in one of them. As a result, the space of maps  $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$  is *first-countable* with respect to the topology of compact convergence, i.e. each point has a countable base of neighborhoods, and a point  $f$  is a limit point of a set  $S$  if and only if there is a sequence of points in  $S$  converging to  $f$ . In a general topological space, however, a limit point of a set  $S$  may not be representable as the limit of a sequence of points from  $S$ . In particular, the space  $L^2(\mathbb{R}^v, \mathbb{R}^{dv})$  is not hemicompact and the space of maps  $f : L^2(\mathbb{R}^v, \mathbb{R}^{dv}) \rightarrow L^2(\mathbb{R}^v, \mathbb{R}^{dv})$  is not first countable with respect to the topology of compact convergence, so that, in particular, we must distinguish between the notions of limit points of the set of convnets and the limits of sequences of convnets. We refer the reader, e.g., to the book [30] for a general discussion of this and other topological questions and in particular to §46 for a discussion of compact convergence.

When defining a limiting map, we would like to require the convnets to increase their resolution  $\frac{1}{\lambda}$  and range  $\Lambda$ . At the same time, we will regard the receptive field and its range parameter  $L_{\text{rf}}$  as arbitrary but fixed (the current common practice in applications is to use small values such as  $L_{\text{rf}} = 1$  regardless of the size of the network; see, e.g., the architecture of residual networks [13] providing state-of-the-art performance on image recognition tasks).

With all these considerations in mind, we introduce the following definition of a limit point of convnets.

**Definition 3.2** With  $V = L^2(\mathbb{R}^v, \mathbb{R}^{dv})$  and  $U = L^2(\mathbb{R}^v, \mathbb{R}^{du})$ , we say that a map  $f : V \rightarrow U$  is a **limit point of basic convnets** if for any  $L_{\text{rf}}$ , any compact set  $K \subset V$ , and any  $\epsilon > 0, \lambda_0 > 0$  and  $\Lambda_0 > 0$  there exists a basic convnet  $\hat{f}$  with the receptive field parameter  $L_{\text{rf}}$ , spacing  $\lambda \leq \lambda_0$  and range  $\Lambda \geq \Lambda_0$  such that  $\sup_{\Phi \in K} \|\hat{f}(\Phi) - f(\Phi)\| < \epsilon$ .

We can state now the main result of this section.

**Theorem 3.1** *A map  $f : V \rightarrow U$  is a limit point of basic convnets if and only if  $f$  is  $\mathbb{R}^v$ -equivariant and continuous in the norm topology.*

Before giving the proof of the theorem, we recall the useful notion of *strong convergence* of linear operators on Hilbert spaces. Namely, if  $A_n$  is a sequence of bounded linear operators on a Hilbert space and  $A$  is another such operator, then we say that the sequence  $A_n$  converges strongly to  $A$  if  $A_n \Phi$  converges to  $A \Phi$  for any vector  $\Phi$  from this Hilbert space. More generally, strong convergence can be defined, by the same reduction, for any family  $\{A_\alpha\}$  of linear operators once the convergence of the family of vectors  $\{A_\alpha \Phi\}$  is specified.

An example of a strongly convergent family is the family of discretizing projectors  $P_\lambda$  defined in (3.6). These projectors converge strongly to the identity as the grid spacing tends to 0:  $P_\lambda \Phi \xrightarrow{\lambda \rightarrow 0} \Phi$ . Another example is the family of projectors  $P_{\lambda, \Lambda}$  projecting  $V$  onto the subspace  $V_{\lambda, \Lambda}$  of discretized and cut-off signals defined in (3.8). It is easy to see that  $P_{\lambda, \Lambda}$  converge strongly to the identity as the spacing tends to 0 and the cutoff is lifted, i.e. as  $\lambda \rightarrow 0$  and  $\Lambda \rightarrow \infty$ . Finally, our representations  $R_\gamma$

defined in (3.5) are strongly continuous in the sense that  $R_{\gamma'}$  converges strongly to  $R_\gamma$  as  $\gamma' \rightarrow \gamma$ .

A useful standard tool in proving strong convergence is the *continuity argument*: if the family  $\{A_\alpha\}$  is uniformly bounded, then the convergence  $A_\alpha \Phi \rightarrow A\Phi$  holds for all vectors  $\Phi$  from the Hilbert space once it holds for a dense subset of vectors. This follows by approximating any  $\Phi$  with  $\Psi$ 's from the dense subset and applying the inequality  $\|A_\alpha \Phi - A\Phi\| \leq \|A_\alpha \Psi - A\Psi\| + (\|A_\alpha\| + \|A\|)\|\Phi - \Psi\|$ . In the sequel, we will consider strong convergence only in the settings where  $A_\alpha$  are orthogonal projectors or norm-preserving operators, so the continuity argument will be applicable.

**Proof of Theorem 3.1 Necessity** (a limit point of basic convnets is  $\mathbb{R}^v$ -equivariant and continuous).

We start by noting that basic convnets  $\hat{f} : V \rightarrow U$ , as defined by Definition 3.1, are continuous in the norm topology, because the initial projection  $P_{\lambda, \Lambda + (T-1)\lambda L_{\text{rf}}}$  is continuous, and the subsequent layers are finite-dimensional transformations that are continuous since they are composed of linear operations and continuous activation functions. The continuity of a limit point of convnets then follows from the continuity of convnets and their uniform convergence on compact sets by a standard argument (see Theorem 46.5 in Munkres [30]).

Let  $f$  denote a limit point of convnets. Let us prove the  $\mathbb{R}^v$ -equivariance of  $f$ , i.e.

$$f(R_\gamma \Phi) = R_\gamma f(\Phi), \quad \gamma \in \mathbb{R}^v, \Phi \in V. \tag{3.13}$$

Let  $D_M = [-M, M]^v \subset \mathbb{R}^v$  with some  $M > 0$ , and  $P_{D_M}$  be the orthogonal projector in  $U$  onto the subspace of signals supported on the set  $D_M$ . Then  $P_{D_M}$  converges strongly to the identity as  $M \rightarrow +\infty$ . Hence, (3.13) will follow if we prove that for any  $M$

$$P_{D_M} f(R_\gamma \Phi) = P_{D_M} R_\gamma f(\Phi). \tag{3.14}$$

Let  $\epsilon > 0$ . Let  $\gamma_\lambda \in (\lambda\mathbb{Z})^v$  be the nearest point to  $\gamma \in \mathbb{R}^v$  on the grid  $(\lambda\mathbb{Z})^v$ . Then, since  $R_{\gamma_\lambda}$  converges strongly to  $R_\gamma$  as  $\lambda \rightarrow 0$ , there exist  $\lambda_0$  such that for any  $\lambda < \lambda_0$

$$\|R_{\gamma_\lambda} f(\Phi) - R_\gamma f(\Phi)\| \leq \epsilon, \tag{3.15}$$

and

$$\|f(R_\gamma \Phi) - f(R_{\gamma_\lambda} \Phi)\| \leq \epsilon, \tag{3.16}$$

where we have also used the already proven continuity of  $f$ .

Observe that the discretization/cutoff projectors  $P_{\lambda, M}$  converge strongly to  $P_{D_M}$  as  $\lambda \rightarrow 0$ , hence we can ensure that for any  $\lambda < \lambda_0$  we also have

$$\begin{aligned} \|P_{D_M} f(R_\gamma \Phi) - P_{\lambda, M} f(R_\gamma \Phi)\| &\leq \epsilon, \\ \|P_{\lambda, M} R_\gamma f(\Phi) - P_{D_M} R_\gamma f(\Phi)\| &\leq \epsilon. \end{aligned} \tag{3.17}$$

Next, observe that basic convnets are partially translationally equivariant by our definition, in the sense that if the cutoff parameter  $\Lambda$  of the convnet is sufficiently large then

$$P_{\lambda, M} \widehat{f}(R_{\gamma_\lambda} \Phi) = P_{\lambda, M} R_{\gamma_\lambda} \widehat{f}(\Phi). \tag{3.18}$$

Indeed, note first that, away from the boundary of the domain  $[-\Lambda, \Lambda]^v$ , all the operations of the convnet (the initial discretizing projection and subsequent convolutional layers) are equivariant with respect to the subgroup  $(\lambda\mathbb{Z})^v \subset \mathbb{R}^v$ . Accordingly, since  $\gamma_\lambda \in (\lambda\mathbb{Z})^v$ , we have  $\widehat{f}(R_{\gamma_\lambda} \Phi)(\lambda\mathbf{k}) = R_{\gamma_\lambda} \widehat{f}(\Phi)(\lambda\mathbf{k})$  for any point  $\lambda\mathbf{k}$  such that both  $\lambda\mathbf{k}$  and  $\lambda\mathbf{k} - \gamma_\lambda$  belong to the convnet output domain  $\lambda Z_{\lfloor \Lambda/\lambda \rfloor}$ . In other words, Eq. (3.18) holds as long as both sets  $\lambda Z_{\lfloor M/\lambda \rfloor}$  and  $\lambda Z_{\lfloor M/\lambda \rfloor} - \gamma_\lambda$  are subsets of  $\lambda Z_{\lfloor \Lambda/\lambda \rfloor}$ . This condition is satisfied if we require that  $\Lambda > \Lambda_0$  with  $\Lambda_0 = M + \|\gamma\|_\infty$ .

Now, take the compact set  $K = \{R_\theta \Phi \mid \theta \in \mathcal{N}\}$ , where  $\mathcal{N} \subset \mathbb{R}^v$  is some compact set including 0 and all points  $\gamma_\lambda$  for  $\lambda < \lambda_0$ . Then, by our definition of a limit point of basic convnets, there is a convnet  $\widehat{f}$  with  $\lambda < \lambda_0$  and  $L > L_0$  such that for all  $\theta \in \mathcal{N}$  (and in particular for  $\theta = 0$  or  $\theta = \gamma_\lambda$ )

$$\|f(R_\theta \Phi) - \widehat{f}(R_\theta \Phi)\| < \epsilon. \tag{3.19}$$

We can now write a bound for the difference of the two sides of (3.14):

$$\begin{aligned} & \|P_{D_M} f(R_\gamma \Phi) - P_{D_M} R_\gamma f(\Phi)\| \\ & \leq \|P_{D_M} f(R_\gamma \Phi) - P_{\lambda, M} f(R_\gamma \Phi)\| + \|P_{\lambda, M} f(R_\gamma \Phi) - P_{\lambda, M} f(R_{\gamma_\lambda} \Phi)\| \\ & \quad + \|P_{\lambda, M} f(R_{\gamma_\lambda} \Phi) - P_{\lambda, M} \widehat{f}(R_{\gamma_\lambda} \Phi)\| + \|P_{\lambda, M} \widehat{f}(R_{\gamma_\lambda} \Phi) - P_{\lambda, M} R_{\gamma_\lambda} \widehat{f}(\Phi)\| \\ & \quad + \|P_{\lambda, M} R_{\gamma_\lambda} \widehat{f}(\Phi) - P_{\lambda, M} R_{\gamma_\lambda} f(\Phi)\| + \|P_{\lambda, M} R_{\gamma_\lambda} f(\Phi) - P_{\lambda, M} R_\gamma f(\Phi)\| \\ & \quad + \|P_{\lambda, M} R_\gamma f(\Phi) - P_{D_M} R_\gamma f(\Phi)\| \\ & \leq \|P_{D_M} f(R_\gamma \Phi) - P_{\lambda, M} f(R_\gamma \Phi)\| + \|f(R_\gamma \Phi) - f(R_{\gamma_\lambda} \Phi)\| \\ & \quad + \|f(R_{\gamma_\lambda} \Phi) - \widehat{f}(R_{\gamma_\lambda} \Phi)\| + \|\widehat{f}(\Phi) - f(\Phi)\| \\ & \quad + \|R_{\gamma_\lambda} f(\Phi) - R_\gamma f(\Phi)\| + \|P_{\lambda, M} R_\gamma f(\Phi) - P_{D_M} R_\gamma f(\Phi)\| \\ & \leq 6\epsilon, \end{aligned}$$

Here in the first step we split the difference into several parts, in the second step we used the identity (3.18) and the fact that  $P_{\lambda, M}, R_{\gamma_\lambda}$  are linear operators with the operator norm 1, and in the third step we applied the inequalities (3.15)–(3.17) and (3.19). Since  $\epsilon$  was arbitrary, we have proved (3.14).

**Sufficiency** (an  $\mathbb{R}^v$ -equivariant and continuous map is a limit point of basic convnets).

We start by proving a key lemma on the approximation capability of basic convnets in the special case when they have the degenerate output range,  $\Lambda = 0$ . In this case, by (3.9), the output space  $W_T = U_{\lambda, 0} \cong \mathbb{R}^{d_U}$ , and the first auxiliary space  $W_1 = V_{\lambda, (T-1)\lambda L_{\text{rf}}} \subset V$ .

**Lemma 3.1** *Let  $\lambda, T$  be fixed and  $\Lambda = 0$ . Then any continuous map  $f : V_{\lambda, (T-1)\lambda L_{\text{rf}}} \rightarrow U_{\lambda, 0}$  can be approximated by basic convnets having spacing  $\lambda$ , depth  $T$ , and range  $\Lambda = 0$ .*

Note that this is essentially a finite-dimensional approximation result, in the sense that the input space  $V_{\lambda, (T-1)\lambda L_{\text{rf}}}$  is finite-dimensional and fixed. The approximation is achieved by choosing sufficiently large feature dimensions  $d_t$  and suitable weights in the intermediate layers.

**Proof** The idea of the proof is to divide the operation of the convnet into two stages. The first stage is implemented by the first  $T - 2$  layers and consists in approximate “contraction” of the input vectors, while the second stage, implemented by the remaining two layers, performs the actual approximation.

The contraction stage is required because the components of the input signal  $\Phi_{\text{in}} \in V_{\lambda, (T-1)\lambda L_{\text{rf}}} \cong L^2(\lambda Z_{(T-1)L_{\text{rf}}}, \mathbb{R}^{d_V})$  are distributed over the large spatial domain  $\lambda Z_{(T-1)L_{\text{rf}}}$ . In this stage we will map the input signal to the spatially localized space  $W_{T-1} \cong L^2(\lambda Z_{L_{\text{rf}}}, \mathbb{R}^{d_{T-1}})$  so as to approximately preserve the information in the signal.

Regarding the second stage, observe that the last two layers of the convnet (starting from  $W_{T-1}$ ) act on signals in  $W_{T-1}$  by an expression analogous to the one-hidden-layer network from the basic universal approximation theorem (Theorem 2.1):

$$\left(\widehat{f}_T \circ \widehat{f}_{T-1}(\Phi)\right)_n = \sum_{k=1}^{d_T} w_{nk}^{(T)} \sigma\left(\sum_{\theta \in Z_{L_{\text{rf}}}} \sum_{m=1}^{d_{T-1}} w_{k\theta m}^{(T-1)} \Phi_{\theta m} + h_k^{(T-1)}\right) + h_n^{(T)}. \tag{3.20}$$

This expression involves all components of  $\Phi \in W_{T-1}$ , and so we can conclude by Theorem 2.1 that by choosing a sufficiently large dimension  $d_T$  and appropriate weights we can approximate an arbitrary continuous map from  $W_{T-1}$  to  $U_{\lambda, 0}$ .

Now, given a continuous map  $f : V_{\lambda, (T-1)\lambda L_{\text{rf}}} \rightarrow U_{\lambda, 0}$ , consider the map  $g = f \circ I \circ P : W_{T-1} \rightarrow U_{\lambda, 0}$ , where  $I$  is some linear isometric map from a subspace  $W'_{T-1} \subset W_{T-1}$  to  $V_{\lambda, (T-1)\lambda L_{\text{rf}}}$ , and  $P$  is the projection in  $W_{T-1}$  to  $W'_{T-1}$ . Such isometric  $I$  exists if  $\dim W_{T-1} \geq \dim V_{\lambda, (T-1)\lambda L_{\text{rf}}}$ , which we can assume w.l.o.g. by choosing sufficiently large  $d_{T-1}$ . Then the map  $g$  is continuous, and the previous argument shows that we can approximate  $g$  using the second stage of the convnet. Therefore, we can also approximate the given map  $f = g \circ I^{-1}$  by the whole convnet if we manage to exactly implement or approximate the isometry  $I^{-1}$  in the contraction stage.

Implementing such an isometry would be straightforward if the first  $T - 2$  layers had no activation function (i.e., if  $\sigma$  were the identity function in the nonlinear layers (3.11)). In this case for all  $t = 2, 3, \dots, T - 1$  we can choose the feature dimensions  $d_t = |Z_{L_{\text{rf}}}| d_{t-1} = (2L_{\text{rf}} + 1)^{v(t-1)} d_V$  and set  $h_n^{(t)} = 0$  and

$$w_{n\theta k}^{(t)} = \begin{cases} 1, & n = \psi_t(\theta, k), \\ 0, & \text{otherwise,} \end{cases}$$

where  $\psi_t$  is some bijection between  $Z_{L_{rf}} \times \{1, \dots, d_t\}$  and  $\{1, \dots, d_{t+1}\}$ . In this way, each component of the network input vector  $\Phi_{in}$  gets copied, layer by layer, to subsequent layers and eventually ends up among the components of the resulting vector in  $W_{T-1}$  (with some repetitions due to multiple possible trajectories of copying).

However, since  $\sigma$  is not an identity, copying needs to be approximated. Consider the first layer,  $\widehat{f}_1$ . For each  $\gamma \in Z_{L_{rf}}$  and each  $s \in \{1, \dots, d_1\}$ , consider the corresponding coordinate map

$$g_{\gamma s} : L^2(\lambda Z_{L_{rf}}, \mathbb{R}^{d_1}) \rightarrow \mathbb{R}, \quad g_{\gamma s} : \Phi \mapsto \Phi_{\gamma s}.$$

By Theorem 2.1, the map  $g_{\gamma s}$  can be approximated with arbitrary accuracy on any compact set in  $L^2(\lambda Z_{L_{rf}}, \mathbb{R}^{d_1})$  by maps of the form

$$\Phi \mapsto \sum_{m=1}^N c_{\gamma sm} \sigma \left( \sum_{\theta \in Z_{L_{rf}}} \sum_{k=1}^{d_1} w_{\gamma sm \theta k} \Phi_{\theta k} + h_{\gamma sm} \right), \tag{3.21}$$

where we may assume without loss of generality that  $N$  is the same for all  $\gamma, s$ . We then set the second feature dimension  $d_2 = N|Z_{L_{rf}}|d_1 = N(2L_{rf} + 1)^v d_V$  and assign the weights  $w_{\gamma sm \theta k}$  and  $h_{\gamma sm}$  in (3.21) to be the weights  $w_{n\theta k}^{(1)}$  and  $h_n^{(1)}$  of the first convnet layer, where the index  $n$  somehow enumerates the triplets  $(\gamma, s, m)$ . Defined in this way, the first convolutional layer  $f_1$  only partly reproduces the copy operation, since this layer does not include the linear weighting corresponding to the external summation over  $m$  in (3.21). However, we can include this weighting into the next layer, since this operation involves only values at the same spatial location  $\gamma \in \mathbb{Z}^v$ , and prepending this operation to the convolutional layer (3.21) does not change the functional form of the layer.

By repeating this argument for the subsequent layers  $t = 2, 3, \dots, T - 2$ , we can make the sequence of the first  $T - 2$  layers to arbitrarily accurately copy all the components of the input vector  $\Phi_{in}$  into a vector  $\Phi \in W_{T-1}$ , up to some additional linear transformations that need to be included in the  $(T - 1)$ 'th layer (again, this is legitimate since prepending a linear operation does not change the functional form of the  $(T - 1)$ 'th layer). Thus, we can approximate  $f = g \circ I^{-1}$  by arranging the first stage of the convnet to approximate  $I^{-1}$  and the second to approximate  $g$ .  $\square$

Returning to the proof of sufficiency, let  $f : V \rightarrow U$  be an  $\mathbb{R}^v$ -equivariant continuous map that we need to approximate with accuracy  $\epsilon$  on a compact set  $K \subset V$  by a convnet with  $\lambda < \lambda_0$  and  $\Lambda > \Lambda_0$ . For any  $\lambda$  and  $\Lambda$ , define the map

$$f_{\lambda, \Lambda} = P_{\lambda, \Lambda} \circ f \circ P_{\lambda}.$$

Observe that we can find  $\lambda < \lambda_0$  and  $\Lambda > \Lambda_0$  such that

$$\sup_{\Phi \in K} \|f_{\lambda, \Lambda}(\Phi) - f(\Phi)\| \leq \frac{\epsilon}{3}. \tag{3.22}$$

Indeed, this can be proved as follows. Denote by  $B_\delta(\Phi)$  the radius- $\delta$  ball centered at  $\Phi$ . By compactness of  $K$  and continuity of  $f$  we can find finitely many signals  $\Phi_n \in V, n = 1, \dots, N$ , and some  $\delta > 0$  so that, first,  $K \subset \cup_n B_{\delta/2}(\Phi_n)$ , and second,

$$\|f(\Phi) - f(\Phi_n)\| \leq \frac{\epsilon}{9}, \quad \Phi \in B_\delta(\Phi_n). \tag{3.23}$$

For any  $\Phi \in K$ , pick  $n$  such that  $\Phi \in B_{\delta/2}(\Phi_n)$ , then

$$\begin{aligned} \|f_{\lambda,\Lambda}(\Phi) - f(\Phi)\| &\leq \|P_{\lambda,\Lambda}f(P_\lambda\Phi) - P_{\lambda,\Lambda}f(\Phi_n)\| \\ &\quad + \|P_{\lambda,\Lambda}f(\Phi_n) - f(\Phi_n)\| + \|f(\Phi_n) - f(\Phi)\| \\ &\leq \|f(P_\lambda\Phi) - f(\Phi_n)\| + \|P_{\lambda,\Lambda}f(\Phi_n) - f(\Phi_n)\| + \frac{\epsilon}{9}. \end{aligned} \tag{3.24}$$

Since  $\Phi \in B_{\delta/2}(\Phi_n)$ , if  $\lambda$  is sufficiently small then  $P_\lambda\Phi \in B_\delta(\Phi_n)$  (by the strong convergence of  $P_\lambda$  to the identity) and hence  $\|f(P_\lambda\Phi) - f(\Phi_n)\| < \frac{\epsilon}{9}$ , again by (3.23). This choice of  $\lambda$  can be made uniformly in  $\Phi \in K$  thanks to the compactness of  $K$ . Also, we can choose sufficiently small  $\lambda$  and then sufficiently large  $\Lambda$  so that  $\|P_{\lambda,\Lambda}f(\Phi_n) - f(\Phi_n)\| < \frac{\epsilon}{9}$ . Using these inequalities in (3.24), we obtain (3.22).

Having thus chosen  $\lambda$  and  $\Lambda$ , observe that, by translation equivariance of  $f$ , the map  $f_{\lambda,\Lambda}$  can be written as

$$f_{\lambda,\Lambda}(\Phi) = \sum_{\gamma \in Z_{[\Lambda/\lambda]}} R_{\lambda\gamma} P_{\lambda,0} f(P_\lambda R_{-\lambda\gamma} \Phi),$$

where  $P_{\lambda,0}$  is the projector  $P_{\lambda,\Lambda}$  in the degenerate case  $\Lambda = 0$ . Consider the map

$$f_{\lambda,\Lambda,T}(\Phi) = \sum_{\gamma \in Z_{[\Lambda/\lambda]}} R_{\lambda\gamma} P_{\lambda,0} f(P_{\lambda,(T-1)\lambda L_{\text{ff}}} R_{-\lambda\gamma} \Phi).$$

Then, by choosing  $T$  sufficiently large, we can ensure that

$$\sup_{\Phi \in K} \|f_{\lambda,\Lambda,T}(\Phi) - f_{\lambda,\Lambda}(\Phi)\| < \frac{\epsilon}{3}. \tag{3.25}$$

Indeed, this can be proved in the same way as (3.22), by using compactness of  $K$ , continuity of  $f$ , finiteness of  $Z_{[\Lambda/\lambda]}$  and the strong convergence  $P_{\lambda,(T-1)\lambda L_{\text{ff}}} R_{-\lambda\gamma} \Phi \xrightarrow{T \rightarrow \infty} P_\lambda R_{-\lambda\gamma} \Phi$ .

Observe that  $f_{\lambda,\Lambda,T}$  can be alternatively written as

$$f_{\lambda,\Lambda,T}(\Phi) = \sum_{\gamma \in Z_{[\Lambda/\lambda]}} R_{\lambda\gamma} f_{\lambda,0,T}(R_{-\lambda\gamma} \Phi), \tag{3.26}$$

where

$$f_{\lambda,0,T}(\Phi) = P_{\lambda,0}f(P_{\lambda,(T-1)\lambda L_{\text{rf}}}\Phi).$$

We can view the map  $f_{\lambda,0,T}$  as a map from  $V_{\lambda,(T-1)\lambda L_{\text{rf}}}$  to  $U_{\lambda,0}$ , which makes Lemma 3.1 applicable to  $f_{\lambda,0,T}$ . Hence, since  $\cup_{\gamma \in Z_{[\Lambda/\lambda]}} R_{-\lambda\gamma}K$  is compact, we can find a convnet  $\widehat{f}_0$  with spacing  $\lambda$ , depth  $T$  and range  $\Lambda = 0$  such that

$$\|\widehat{f}_0(\Phi) - f_{\lambda,0,T}(\Phi)\| < \frac{\epsilon}{3|Z_{[\Lambda/\lambda]}|}, \quad \Phi \in \cup_{\gamma \in Z_{[\Lambda/\lambda]}} R_{-\lambda\gamma}K. \tag{3.27}$$

Consider the convnet  $\widehat{f}_\Lambda$  different from  $\widehat{f}_0$  only by the range parameter  $\Lambda$ ; such a convnet can be written in terms of  $\widehat{f}_0$  in the same way as  $f_{\lambda,\Lambda,T}$  is written in terms of  $f_{\lambda,0,T}$ :

$$\widehat{f}_\Lambda(\Phi) = \sum_{\gamma \in Z_{[\Lambda/\lambda]}} R_{\lambda\gamma} \widehat{f}_0(R_{-\lambda\gamma}\Phi). \tag{3.28}$$

Combining (3.26), (3.27) and (3.28), we obtain

$$\sup_{\Phi \in K} \|\widehat{f}_\Lambda(\Phi) - f_{\lambda,\Lambda,T}(\Phi)\| < \frac{\epsilon}{3}.$$

Combining this bound with bounds (3.22) and (3.25), we obtain the desired bound

$$\sup_{\Phi \in K} \|\widehat{f}_\Lambda(\Phi) - f(\Phi)\| < \epsilon.$$

□

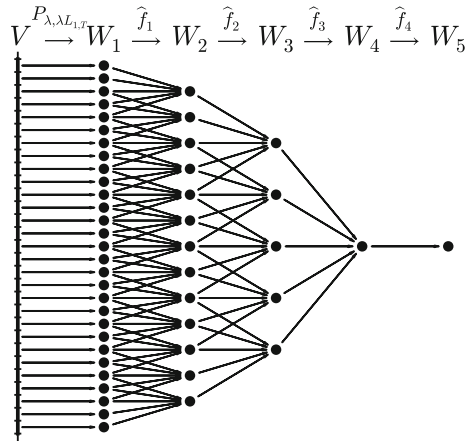
Theorem 3.1 suggests that our definition of limit points of basic convnets provides a reasonable rigorous framework for the analysis of convergence and invariance properties of convnet-like models in the limit of continual and infinitely extended signals. We will use these definition and theorem as templates when considering convnets with pooling in the next subsection and charge-conserving convnets in Sect. 4.

### 3.3 Convnets with Pooling

As already mentioned, pooling erodes the equivariance of models with respect to translations. Therefore, we will consider convnets with pooling as universal approximators without assuming the approximated maps to be translationally invariant. Also, rather than considering  $L^2(\mathbb{R}^v, \mathbb{R}^{d_v})$ -valued maps, we will be interested in approximating simply  $\mathbb{R}$ -valued maps, i.e., those of the form  $f : V \rightarrow \mathbb{R}$ , where, as in Sect. 3.2,  $V = L^2(\mathbb{R}^v, \mathbb{R}^{d_v})$  (Fig. 2).

While the most popular kind of pooling in applications seems to be max-pooling, we will only consider pooling by decimation (i.e., grid downsampling), which appears to be about as efficient in practice (see Springenberg et al. [42]). Compared to basic

**Fig. 2** A one-dimensional ( $v = 1$ ) convnet with downsampling having stride  $s = 2$  and the receptive field parameter  $L_{\text{rf}} = 2$



convnets of Sect. 3.2, convnets with downsampling then have a new parameter, *stride*, that we denote by  $s$ . The stride can take values  $s = 1, 2, \dots$  and determines the geometry scaling when passing information to the next convnet layer: if the current layer operates on a grid  $(\lambda\mathbb{Z})^v$ , then the next layer will operate on the subgrid  $(s\lambda\mathbb{Z})^v$ . Accordingly, the current layer only needs to perform the operations having outputs located in this subgrid. We will assume  $s$  to be fixed and to be the same for all layers. Moreover, we assume that

$$s \leq 2L_{\text{rf}} + 1, \tag{3.29}$$

i.e., the stride is not larger than the size of the receptive field: this ensures that information from each node of the current grid can reach the next layer.

Like the basic convnet of Sect. 3.2, a convnet with downsampling can be written as a chain:

$$\hat{f} : V \xrightarrow{P_{\lambda, \lambda L_{1,T}}} V_{\lambda, \lambda L_{1,T}} (\equiv W_1) \xrightarrow{\hat{f}_1} W_2 \xrightarrow{\hat{f}_2} \dots \xrightarrow{\hat{f}_T} W_{T+1} (\cong \mathbb{R}). \tag{3.30}$$

Here the space  $V_{\lambda, \lambda L_{1,T}}$  is defined as in (3.8) (with  $\Lambda = \lambda L_{1,T}$ ) and  $P_{\lambda, \lambda L_{1,T}}$  is the orthogonal projector to this subspace. The intermediate spaces are defined by

$$W_t = L^2(s^{t-1}\lambda\mathbb{Z}_{L_{1,T}}, \mathbb{R}^{d_t}).$$

The range parameters  $L_{t,T}$  are given by

$$L_{t,T} = \begin{cases} L_{\text{rf}}(1 + s + s^2 + \dots + s^{T-t-1}), & t < T, \\ 0, & t = T, T + 1. \end{cases}$$

This choice of  $L_{t,T}$  is equivalent to the identities

$$L_{t,T} = sL_{t+1,T} + L_{\text{rf}}, \quad t = 1, \dots, T - 1,$$



expressing the domain transformation under downsampling.

The feature dimensions  $d_t$  can again take any values, aside from the fixed values  $d_1 = d_V$  and  $d_{T+1} = 1$ .

As the present convnet model is  $\mathbb{R}$ -valued, in contrast to the basic convnet of Sect. 3.2, it does not have a separate output cutoff parameter  $\Lambda$  (we essentially have  $\Lambda = 0$  now). The geometry of the input domain  $\lambda Z_{L_{1,T}}$  is fully determined by stride  $s$ , the receptive field parameter  $L_{\text{rf}}$ , grid spacing  $\lambda$ , and depth  $T$ . Thus, the architecture of the model is fully specified by these parameters and feature dimensions  $d_2, \dots, d_T$ .

The layer operation formulas differ from the formulas (3.11), (3.3) by the inclusion of downsampling:

$$\widehat{f}_t(\Phi)_{\gamma n} = \sigma \left( \sum_{\theta \in Z_{L_{\text{rf}}}} \sum_{k=1}^{d_t} w_{n\theta k}^{(t)} \Phi_{s\gamma+\theta, k} + h_n^{(t)} \right), \quad \gamma \in Z_{L_{t+1}}, n = 1, \dots, d_{t+1}, \quad t \leq T, \tag{3.31}$$

$$\widehat{f}_{T+1}(\Phi) = \sum_{k=1}^{d_T} w_{nk}^{(T)} \Phi_k + h_n^{(T)}. \tag{3.32}$$

Summarizing, we define convnets with downsampling as follows.

**Definition 3.3** A **convnet with downsampling** is a map  $\widehat{f} : V \rightarrow \mathbb{R}$  defined by (3.30), (3.31), (3.32), and characterized by parameters  $s, \lambda, L_{\text{rf}}, T, d_1, \dots, d_T$  and coefficients  $w_{n\theta k}^{(t)}$  and  $h_n^{(t)}$ .

Next, we give a definition of a limit point of convnets with downsampling analogous to Definition 3.2 for basic convnets. In this definition, we require that the input domain grow in resolution  $\frac{1}{\lambda}$  and in the spatial range  $\lambda L_{1,T}$ , while the stride and receptive field are fixed.

**Definition 3.4** With  $V = L^2(\mathbb{R}^v, \mathbb{R}^{d_V})$ , we say that a map  $f : V \rightarrow \mathbb{R}$  is a **limit point of convnets with downsampling** if for any  $s$  and  $L_{\text{rf}}$  subject to Eq. (3.29), any compact set  $K \in V$ , any  $\epsilon > 0, \lambda_0 > 0$  and  $\Lambda_0 > 0$  there exists a convnet with downsampling  $\widehat{f}$  with stride  $s$ , receptive field parameter  $L_{\text{rf}}$ , depth  $T$ , and spacing  $\lambda \leq \lambda_0$  such that  $\lambda L_{1,T} \geq \Lambda_0$  and  $\sup_{\Phi \in K} \|\widehat{f}(\Phi) - f(\Phi)\| < \epsilon$ .

The analog of Theorem 3.1 then reads:

**Theorem 3.2** A map  $f : V \rightarrow \mathbb{R}$  is a limit point of convnets with downsampling if and only if  $f$  is continuous in the norm topology.

**Proof** The proof is completely analogous to, and in fact simpler than, the proof of Theorem 3.1, so we only sketch it.

The necessity only involves the claim of continuity and follows again by a basic topological argument.

In the proof of sufficiency, an analog of Lemma 3.1 holds for convnets with downsampling, since, thanks to the constraint (3.29) on the stride, all points of the input domain  $\lambda Z_{L_1}$  are connected by the network architecture to the output (though there

are fewer connections now due to pooling), so that our construction of approximate copy operations remains valid.

To approximate  $f : V \rightarrow \mathbb{R}$  on a compact  $K$ , first approximate it by a map  $f \circ P_{\lambda, Z_{L_1, T}}$  with a sufficiently small  $\lambda$  and large  $T$ , then use the lemma to approximate  $f \circ P_{\lambda, Z_{L_1, T}}$  by a convnet. □

### 4 Charge-Conserving Convnets

The goal of the present section is to describe a complete convnet-like model for approximating arbitrary continuous maps equivariant with respect to rigid planar motions. A rigid motion of  $\mathbb{R}^v$  is an affine transformation preserving the distances and the orientation in  $\mathbb{R}^v$ . The group  $SE(v)$  of all such motions can be described as a *semidirect product* of the translation group  $\mathbb{R}^v$  with the special orthogonal group  $SO(v)$ :

$$SE(v) = \mathbb{R}^v \rtimes SO(v).$$

An element of  $SE(v)$  can be represented as a pair  $(\gamma, \theta)$  with  $\gamma \in \mathbb{R}^v$  and  $\theta \in SO(v)$ . The group operations are given by

$$\begin{aligned} (\gamma_1, \theta_1)(\gamma_2, \theta_2) &= (\gamma_1 + \theta_1\gamma_2, \theta_1\theta_2), \\ (\gamma, \theta)^{-1} &= (-\theta^{-1}\gamma, \theta^{-1}). \end{aligned}$$

The group  $SE(v)$  acts on  $\mathbb{R}^v$  by

$$\mathcal{A}_{(\gamma, \theta)}\mathbf{x} = \gamma + \theta\mathbf{x}.$$

It is easy to see that this action is compatible with the group operation, i.e.  $\mathcal{A}_{(0, 1)} = \text{Id}$  and  $\mathcal{A}_{(\gamma_1, \theta_1)}\mathcal{A}_{(\gamma_2, \theta_2)} = \mathcal{A}_{(\gamma_1, \theta_1)(\gamma_2, \theta_2)}$  (implying, in particular,  $\mathcal{A}_{(\gamma, \theta)}^{-1} = \mathcal{A}_{(\gamma, \theta)^{-1}}$ ).

As in Sect. 3.2, consider the space  $V = L^2(\mathbb{R}^v, \mathbb{R}^{d_V})$ . We can view this space as a  $SE(v)$ -module with the representation canonically associated with the action  $\mathcal{A}$ :

$$R_{(\gamma, \theta)}\Phi(\mathbf{x}) = \Phi(\mathcal{A}_{(\gamma, \theta)^{-1}}\mathbf{x}), \tag{4.1}$$

where  $\Phi : \mathbb{R}^v \rightarrow \mathbb{R}^{d_V}$  and  $\mathbf{x} \in \mathbb{R}^v$ . We define in the same manner the module  $U = L^2(\mathbb{R}^v, \mathbb{R}^{d_U})$ . In the remainder of the paper we will be interested in approximating *continuous and  $SE(v)$ -equivariant* maps  $f : V \rightarrow U$ . Let us first give some examples of such maps.

**Linear maps.** Assume for simplicity that  $d_V = d_U = 1$  and consider a *linear*  $SE(v)$ -equivariant map  $f : L^2(\mathbb{R}^v) \rightarrow L^2(\mathbb{R}^v)$ . Such a map can be written as a convolution  $f(\Phi) = \Phi * \Psi_f$ , where  $\Psi_f$  is a radial signal,  $\Psi_f(\mathbf{x}) = \tilde{\Psi}_f(|\mathbf{x}|)$ . In general,  $\Psi_f$  should be understood in a distributional sense.

By applying Fourier transform  $\mathcal{F}$ , the map  $f$  can be equivalently described in the Fourier dual space as pointwise multiplication of the given signal by  $const\mathcal{F}\Psi_f$

(with the constant depending on the choice of the coefficient in the Fourier transform), so  $f$  is  $SE(\nu)$ -equivariant and continuous if and only if  $\mathcal{F}\Psi_f$  is a radial function belonging to  $L^\infty(\mathbb{R}^\nu)$ . Note that in this argument we have tacitly complexified the space  $L^2(\mathbb{R}^\nu, \mathbb{R})$  into  $L^2(\mathbb{R}^\nu, \mathbb{C})$ . The condition that  $f$  preserves real-valuedness of the signal  $\Phi$  translates into  $\overline{\mathcal{F}\Psi_f(\mathbf{x})} = \mathcal{F}\Psi_f(-\mathbf{x})$ , where the bar denotes complex conjugation.

Note that linear  $SE(\nu)$ -equivariant differential operators, such as the Laplacian  $\Delta$ , are not included in our class of maps, since they are not even defined on the whole space  $V = L^2(\mathbb{R}^\nu)$ . However, if we consider a smoothed version of the Laplacian given by  $f : \Phi \mapsto \Delta(\Phi * g_\epsilon)$ , where  $g_\epsilon$  is the variance- $\epsilon$  Gaussian kernel, then this map will be well-defined on the whole  $V$ , norm-continuous and  $SE(\nu)$ -equivariant.

**Pointwise maps.** Consider a *pointwise* map  $f : V \rightarrow U$  defined by  $f(\Phi)(\mathbf{x}) = f_0(\Phi(\mathbf{x}))$ , where  $f_0 : \mathbb{R}^{d_V} \rightarrow \mathbb{R}^{d_U}$  is some map. In this case  $f$  is  $SE(\nu)$ -equivariant. Note that if  $f_0(0) \neq 0$ , then  $f$  is not well-defined on  $V = L^2(\mathbb{R}^\nu, \mathbb{R}^{d_V})$ , since  $f(\Phi) \notin L^2(\mathbb{R}^\nu, \mathbb{R}^{d_U})$  for the trivial signal  $\Phi(\mathbf{x}) \equiv 0$ . An easy-to-check sufficient condition for  $f$  to be well-defined and continuous on the whole  $V$  is that  $f_0(0) = 0$  and  $f_0$  be globally Lipschitz (i.e.,  $|f_0(\mathbf{x}) - f_0(\mathbf{y})| \leq c|\mathbf{x} - \mathbf{y}|$  for all  $\mathbf{x}, \mathbf{y} \in \mathbb{R}^\nu$  and some  $c < \infty$ ).

Our goal in this section is to describe a finite computational model that would be a universal approximator for all continuous and  $SE(\nu)$ -equivariant maps  $f : V \rightarrow U$ . Following the strategy of Sect. 3.2, we aim to define limit points of such finite models and then prove that the limit points are exactly the continuous and  $SE(\nu)$ -equivariant maps.

We focus on approximating  $L^2(\mathbb{R}^\nu, \mathbb{R}^{d_U})$ -valued  $SE(\nu)$ -equivariant maps rather than  $\mathbb{R}^{d_U}$ -valued  $SE(\nu)$ -invariant maps because, as discussed in Sect. 3, we find it hard to reconcile the  $SE(\nu)$ -invariance with pooling.

Note that, as in the previous sections, there is a straightforward symmetrization-based approach to constructing universal  $SE(\nu)$ -equivariant models. In particular, the group  $SE(\nu)$  extends the group of translations  $\mathbb{R}^\nu$  by the compact group  $SO(\nu)$ , and we can construct  $SE(\nu)$ -equivariant maps simply by symmetrizing  $\mathbb{R}^\nu$ -equivariant maps over  $SO(\nu)$ , as in Proposition 2.2.

**Proposition 4.1** *If a map  $f_{\mathbb{R}^\nu} : V \rightarrow U$  is continuous and  $\mathbb{R}^\nu$ -equivariant, then the map  $f_{SE(\nu)} : V \rightarrow U$  defined by*

$$f_{SE(\nu)}(\Phi) = \int_{SO(\nu)} R_{(0,\theta)^{-1}} f_{\mathbb{R}^\nu}(R_{(0,\theta)} \Phi) d\theta$$

*is continuous and  $SE(\nu)$ -equivariant.*

**Proof** The continuity of  $f_{SE(\nu)}$  follows by elementary arguments using the continuity of  $f_{\mathbb{R}^\nu} : V \rightarrow U$ , uniform boundedness of the operators  $R_{(0,\theta)}$ , and compactness of  $SO(\nu)$ . The  $SE(\nu)$ -equivariance follows since for any  $(\gamma, \theta') \in SE(\nu)$  and  $\Phi \in V$

$$f_{SE(\nu)}(R_{(\gamma,\theta')} \Phi) = \int_{SO(\nu)} R_{(0,\theta)^{-1}} f_{\mathbb{R}^\nu}(R_{(0,\theta)} R_{(\gamma,\theta')} \Phi) d\theta$$

$$\begin{aligned}
&= \int_{\text{SO}(v)} R_{(0,\theta)^{-1}} f_{\mathbb{R}^v}(R_{(\theta\gamma,1)} R_{(0,\theta\theta')} \Phi) d\theta \\
&= \int_{\text{SO}(v)} R_{(0,\theta)^{-1}} R_{(\theta\gamma,1)} f_{\mathbb{R}^v}(R_{(0,\theta\theta')} \Phi) d\theta \\
&= \int_{\text{SO}(v)} R_{(\gamma,\theta')} R_{(0,\theta\theta')^{-1}} f_{\mathbb{R}^v}(R_{(0,\theta\theta')} \Phi) d\theta \\
&= R_{(\gamma,\theta')} f_{\text{SE}(v)}(\Phi).
\end{aligned}$$

□

This proposition implies, in particular, that  $\text{SO}(v)$ -symmetrizations of merely  $\mathbb{R}^v$ -equivariant basic convnets considered in Sect. 3.2 can serve as universal  $\text{SE}(v)$ -equivariant approximators. However, like in the previous sections, we will be instead interested in an intrinsically  $\text{SE}(v)$ -equivariant network construction not involving explicit symmetrization of the approximation over the group  $\text{SO}(v)$ . In particular, our approximators will not use rotated grids.

Our construction relies heavily on the representation theory of the group  $\text{SO}(v)$ , and in the present paper we restrict ourselves to the case  $v = 2$ , in which the group  $\text{SO}(v)$  is abelian and the representation theory is much easier than in the general case.

Section 4.1 contains preliminary considerations suggesting the network construction appropriate for our purpose. The formal detailed description of the model is given in Sect. 4.2. In Sect. 4.3 we formulate and prove the main result of the section, the  $\text{SE}(2)$ -equivariant universal approximation property of the model.

## 4.1 Preliminary Considerations

In this section we explain the idea behind our construction of the universal  $\text{SE}(2)$ -equivariant convnet (to be formulated precisely in Sect. 4.2). We start by showing in Sect. 4.1.1 that a  $\text{SE}(2)$ -equivariant map  $f : V \rightarrow U$  can be described using a  $\text{SO}(2)$ -invariant map  $f_{\text{loc}} : V \rightarrow \mathbb{R}^{d_v}$ . Then, relying on this observation, in Sect. 4.1.2 we show that, heuristically,  $f$  can be reconstructed by first equivariantly extracting local “features” from the original signal using equivariant differentiation, and then transforming these features using a  $\text{SO}(2)$ -invariant pointwise map. In Sect. 4.1.3 we describe discretized differential operators and smoothing operators that we require in order to formulate our model as a finite computation model with sufficient regularity. Finally, in Sect. 4.1.4 we consider polynomial approximations on  $\text{SO}(2)$ -modules.

### 4.1.1 Pointwise Characterization of $\text{SE}(v)$ -Equivariant Maps

In this subsection we show that, roughly speaking,  $\text{SE}(v)$ -equivariant maps  $f : V \rightarrow U$  can be described in terms of  $\text{SO}(v)$ -invariant maps  $f : V \rightarrow \mathbb{R}^v$  obtained by observing the output signal at a fixed position.

(The proposition below has one technical subtlety: we consider signal values  $\Phi(0)$  at a particular point  $\mathbf{x} = 0$  for generic signals  $\Phi$  from the space  $L^2(\mathbb{R}^v, \mathbb{R}^{d_U})$ . Elements of this spaces are defined as equivalence classes of signals that can differ on sets of zero

Lebesgue measure, so, strictly speaking,  $\Phi(0)$  is not well-defined. We can circumvent this difficulty by fixing a particular canonical representative of the equivalence class, say

$$\Phi_{\text{canon}}(\mathbf{x}) = \begin{cases} \lim_{\epsilon \rightarrow 0} \frac{1}{|B_\epsilon(\mathbf{x})|} \int_{B_\epsilon(\mathbf{x})} \Phi(\mathbf{y}) \, d\mathbf{y}, & \text{if the limit exists,} \\ 0, & \text{otherwise.} \end{cases}$$

Lebesgue’s differentiation theorem ensures that the limit exists and agrees with  $\Phi$  almost everywhere, so that  $\Phi_{\text{canon}}$  is indeed a representative of the equivalence class. This choice of the representative is clearly  $\text{SE}(v)$ -equivariant. In the proposition below, the signal value at  $\mathbf{x} = 0$  can be understood as the value of such a canonical representative.<sup>1)</sup>

**Proposition 4.2** *Let  $f : L^2(\mathbb{R}^v, \mathbb{R}^{dv}) \rightarrow L^2(\mathbb{R}^v, \mathbb{R}^{dv})$  be a  $\mathbb{R}^v$ -equivariant map. Then  $f$  is  $\text{SE}(v)$ -equivariant if and only if  $f(R_{(0,\theta)}\Phi)(0) = f(\Phi)(0)$  for all  $\theta \in \text{SO}(v)$  and  $\Phi \in V$ .*

**Proof** One direction of the statement is obvious: if  $f$  is  $\text{SE}(v)$ -equivariant, then  $f(R_{(0,\theta)}\Phi)(0) = R_{(0,\theta)}f(\Phi)(0) = f(\Phi)(\mathcal{A}_{(0,\theta^{-1})}0) = f(\Phi)(0)$ .

Let us prove the opposite implication, i.e. that  $f(R_{(0,\theta)}\Phi)(0) \equiv f(\Phi)(0)$  implies the  $\text{SE}(v)$ -equivariance. We need to show that for all  $(\gamma, \theta) \in \text{SE}(v)$ ,  $\Phi \in V$  and  $\mathbf{x} \in \mathbb{R}^v$  we have

$$f(R_{(\gamma,\theta)}\Phi)(\mathbf{x}) = R_{(\gamma,\theta)}f(\Phi)(\mathbf{x}).$$

Indeed,

$$\begin{aligned} f(R_{(\gamma,\theta)}\Phi)(\mathbf{x}) &= R_{(-\mathbf{x},1)}f(R_{(\gamma,\theta)}\Phi)(0) \\ &= f(R_{(-\mathbf{x},1)}R_{(\gamma,\theta)}\Phi)(0) \\ &= f(R_{(0,\theta)}R_{(\theta^{-1}(\gamma-\mathbf{x}),1)}\Phi)(0) \\ &= f(R_{(\theta^{-1}(\gamma-\mathbf{x}),1)}\Phi)(0) \\ &= R_{(\theta^{-1}(\gamma-\mathbf{x}),1)}f(\Phi)(0) \\ &= R_{(\mathbf{x},\theta)}R_{(\theta^{-1}(\gamma-\mathbf{x}),1)}f(\Phi)(\mathcal{A}_{(\mathbf{x},\theta)}0) \\ &= R_{(\gamma,\theta)}f(\Phi)(\mathbf{x}), \end{aligned}$$

where we used definition (4.1) (steps 1 and 6), the  $\mathbb{R}^v$ -equivariance of  $f$  (steps 2 and 5), and the hypothesis of the lemma (step 4). □

<sup>1</sup> Another approach to ensure a well-defined value  $\Phi(\mathbf{x})$  is to work with shift-invariant reproducing kernel Hilbert spaces (RKHS) instead of  $L^2$  spaces. Definition of RKHS requires the signal evaluation  $\Phi \mapsto \Phi(\mathbf{x})$  to be continuous in  $\Phi$  and in particular well-defined. An example of a shift-invariant RKHS is the space of band-limited signals with a particular bandwidth. We thank the anonymous reviewer for pointing out this approach.

Now, if  $f : V \rightarrow U$  is an  $SE(\nu)$ -equivariant map, then we can define the  $SO(\nu)$ -invariant map  $f_{loc} : V \rightarrow \mathbb{R}^{d_U}$  by

$$f_{loc}(\Phi) = f(\Phi)(0). \tag{4.2}$$

Conversely, suppose that  $f_{loc} : V \rightarrow \mathbb{R}^{d_U}$  is an  $SO(\nu)$ -invariant map. Consider the map  $f : V \rightarrow \{\Psi : \mathbb{R}^\nu \rightarrow \mathbb{R}^{d_U}\}$  defined by

$$f(\Phi)(\mathbf{x}) := f_{loc}(R_{(-\mathbf{x},1)}\Phi). \tag{4.3}$$

In general,  $f(\Phi)$  need not be in  $L^2(\mathbb{R}^\nu, \mathbb{R}^{d_U})$ . Suppose, however, that this is the case for all  $\Phi \in V$ . Then  $f$  is clearly  $\mathbb{R}^\nu$ -equivariant and, moreover,  $SE(\nu)$ -equivariant, by the above proposition.

Thus, under some additional regularity assumption, the task of reconstructing  $SE(\nu)$ -equivariant maps  $f : V \rightarrow U$  is equivalent to the task of reconstructing  $SO(\nu)$ -invariant maps  $f_{loc} : V \rightarrow \mathbb{R}^{d_U}$ .

From this point on, we set  $\nu = 2$ .

### 4.1.2 Equivariant Differentiation

It is convenient to describe rigid motions of  $\mathbb{R}^2$  by identifying this two-dimensional real space with the one-dimensional complex space  $\mathbb{C}$ . Then an element of  $SE(2)$  can be written as  $(\gamma, \theta) = (x + iy, e^{i\phi})$  with some  $x, y \in \mathbb{R}$  and  $\phi \in [0, 2\pi)$ . The action of  $SE(2)$  on  $\mathbb{R}^2 \cong \mathbb{C}$  can be written as

$$A_{(x+iy, e^{i\phi})}z = x + iy + e^{i\phi}z, \quad z \in \mathbb{C}.$$

Using analogous notation  $R_{(x+iy, e^{i\phi})}$  for the canonically associated representation of  $SE(2)$  in  $V$  defined in (4.1), consider the generators of this representation:

$$J_x = i \lim_{\delta x \rightarrow 0} \frac{R_{(\delta x, 1)} - 1}{\delta x}, \quad J_y = i \lim_{\delta y \rightarrow 0} \frac{R_{(i\delta y, 1)} - 1}{\delta y}, \quad J_\phi = i \lim_{\delta\phi \rightarrow 0} \frac{R_{(0, e^{i\delta\phi})} - 1}{\delta\phi}.$$

The generators can be explicitly written as

$$J_x = -i\partial_x, \quad J_y = -i\partial_y, \quad J_\phi = -i\partial_\phi = -i(x\partial_y - y\partial_x)$$

and obey the commutation relations

$$[J_x, J_y] = 0, \quad [J_x, J_\phi] = -iJ_y, \quad [J_y, J_\phi] = iJ_x. \tag{4.4}$$

We are interested in local transformations of signals  $\Phi \in V$ , so it is natural to consider the action of differential operators on the signals. We would like, however, to ensure

the equivariance of this action. This can be done as follows. Consider the first-order operators

$$\partial_z = \frac{1}{2}(\partial_x - i\partial_y), \quad \partial_{\bar{z}} = \frac{1}{2}(\partial_x + i\partial_y).$$

These operators commute with  $J_x, J_y$ , and have the following commutation relations with  $J_\phi$ :

$$[\partial_z, J_\phi] = \partial_z, \quad [\partial_{\bar{z}}, J_\phi] = -\partial_{\bar{z}}$$

or, equivalently,

$$\partial_z J_\phi = (J_\phi + 1)\partial_z, \quad \partial_{\bar{z}} J_\phi = (J_\phi - 1)\partial_{\bar{z}}. \tag{4.5}$$

Let us define, for any  $\mu \in \mathbb{Z}$ ,

$$J_\phi^{(\mu)} = J_\phi + \mu = \mu - i\partial_\phi.$$

Then the triple  $(J_x, J_y, J_z^{(\mu)})$  obeys the same commutation relations (4.4), i.e., constitutes another representation of the Lie algebra of the group  $SE(2)$ . The corresponding representation of the group differs from the original representation (4.1) by the extra phase factor:

$$R_{(x+iy, e^{i\phi})}^{(\mu)} \Phi(\mathbf{x}) = e^{-i\mu\phi} \Phi(\mathcal{A}_{(x+iy, e^{i\phi})}^{-1}\mathbf{x}). \tag{4.6}$$

The identities (4.5) imply  $\partial_z J_\phi^{(\mu)} = J_\phi^{(\mu+1)}\partial_z$  and  $\partial_{\bar{z}} J_\phi^{(\mu)} = J_\phi^{(\mu-1)}\partial_{\bar{z}}$ . Since the operators  $\partial_z, \partial_{\bar{z}}$  also commute with  $J_x, J_y$ , we see that the operators  $\partial_z, \partial_{\bar{z}}$  can serve as *ladder operators* equivariantly mapping

$$\partial_z : V_\mu \rightarrow V_{\mu+1}, \quad \partial_{\bar{z}} : V_\mu \rightarrow V_{\mu-1}, \tag{4.7}$$

where  $V_\mu$  is the space  $L^2(\mathbb{R}^2, \mathbb{R}^{dv})$  equipped with the representation (4.6). Thus, we can equivariantly differentiate signals as long as we appropriately switch the representation. In the sequel, we will for brevity refer to the parameter  $\mu$  characterizing the representation as its *global charge*.

It is convenient to also consider another kind of charge, associated with angular dependence of the signal with respect to rotations about fixed points; let us call it *local charge*  $\eta$  in contrast to the above global charge  $\mu$ . Namely, for any fixed  $\mathbf{x}_0 \in \mathbb{R}^2$ , decompose the module  $V_\mu$  as

$$V_\mu = \bigoplus_{\eta \in \mathbb{Z}} V_{\mu, \eta}^{(\mathbf{x}_0)}, \tag{4.8}$$

where

$$V_{\mu,\eta}^{(\mathbf{x}_0)} = R_{(\mathbf{x}_0,1)} V_{\mu,\eta}^{(0)}, \tag{4.9}$$

and

$$V_{\mu,\eta}^{(0)} = \{\Phi \in V_\mu \mid \Phi(\mathcal{A}_{(0,e^{i\phi})^{-1}} \mathbf{x}) = e^{-i\eta\phi} \Phi(\mathbf{x}) \ \forall \phi\}. \tag{4.10}$$

Writing  $\mathbf{x}_0 = (x_0, y_0)$ , we can characterize  $V_{\mu,\eta}^{(\mathbf{x}_0)}$  as the eigenspace of the operator

$$J_\phi^{(\mathbf{x}_0)} := R_{(\mathbf{x}_0,1)} J_\phi R_{(\mathbf{x}_0,1)}^{-1} = -i(x - x_0)\partial_y + i(y - y_0)\partial_x$$

corresponding to the eigenvalue  $\eta$ . The operator  $J_\phi^{(\mathbf{x}_0)}$  has the same commutation relations with  $\partial_z, \partial_{\bar{z}}$  as  $J_\phi$ :

$$[\partial_z, J_\phi^{(\mathbf{x}_0)}] = \partial_z, \quad [\partial_{\bar{z}}, J_\phi^{(\mathbf{x}_0)}] = -\partial_{\bar{z}}.$$

We can then describe the structure of equivariant maps (4.7) with respect to decomposition (4.8) as follows: for any  $\mathbf{x}_0$ , the decrease or increase of the global charge by the respective ladder operator is compensated by the opposite effect of this operator on the local charge, i.e.  $\partial_z$  maps  $V_{\mu,\eta}^{(\mathbf{x}_0)}$  to  $V_{\mu+1,\eta-1}^{(\mathbf{x}_0)}$  while  $\partial_{\bar{z}}$  maps  $V_{\mu,\eta}^{(\mathbf{x}_0)}$  to  $V_{\mu-1,\eta+1}^{(\mathbf{x}_0)}$ :

$$\partial_z V_{\mu,\eta}^{(\mathbf{x}_0)} \rightarrow V_{\mu+1,\eta-1}^{(\mathbf{x}_0)}, \quad \partial_{\bar{z}} : V_{\mu,\eta}^{(\mathbf{x}_0)} \rightarrow V_{\mu-1,\eta+1}^{(\mathbf{x}_0)}. \tag{4.11}$$

We interpret these identities as *conservation of the total charge*,  $\mu + \eta$ . We remark that there is some similarity between our total charge and the total angular momentum in quantum mechanics; the total angular momentum there consists of the spin component and the orbital component that are analogous to our global and local charge, respectively.

Now we give a heuristic argument showing how to express an arbitrary equivariant map  $f : V \rightarrow U$  using our equivariant differentiation. As discussed in the previous subsection, the task of expressing  $f$  reduces to expressing  $f_{\text{loc}}$  using formulas (4.2), (4.3). Let a signal  $\Phi$  be analytic as a function of the real variables  $x, y$ , then it can be Taylor expanded as

$$\Phi = \sum_{a,b=0}^{\infty} \frac{1}{a!b!} \partial_z^a \partial_{\bar{z}}^b \Phi(0) \Phi_{a,b}, \tag{4.12}$$

with the basis signals  $\Phi_{a,b}$  given by

$$\Phi_{a,b}(z) = z^a \bar{z}^b.$$



The signal  $\Phi$  is fully determined by the coefficients  $\partial_z^a \partial_{\bar{z}}^b \Phi(0)$ , so the map  $f_{loc}$  can be expressed as a function of these coefficients:

$$f_{loc}(\Phi) = \tilde{f}_{loc}((\partial_z^a \partial_{\bar{z}}^b \Phi(0))_{a,b=0}^\infty). \tag{4.13}$$

At  $\mathbf{x}_0 = 0$ , the signals  $\Phi_{a,b}$  have local charge  $\eta = a - b$ , and, if viewed as elements of  $V_{\mu=0}$ , transform under rotations by

$$R_{(0,e^{i\phi})} \Phi_{a,b} = e^{-i(a-b)\phi} \Phi_{a,b}.$$

Accordingly, if we write  $\Phi$  in the form  $\Phi = \sum_{a,b} c_{a,b} \Phi_{a,b}$ , then

$$R_{(0,e^{i\phi})} \Phi = \sum_{a,b} e^{-i(a-b)\phi} c_{a,b} \Phi_{a,b}.$$

It follows that the SO(2)-invariance of  $f_{loc}$  is equivalent to  $\tilde{f}_{loc}$  being invariant with respect to simultaneous multiplication of the arguments by the factors  $e^{-i(a-b)\phi}$ :

$$\tilde{f}_{loc}((e^{-i(a-b)\phi} c_{a,b})_{a,b=0}^\infty) = \tilde{f}_{loc}((c_{a,b})_{a,b=0}^\infty) \quad \forall \phi.$$

Having determined the invariant map  $\tilde{f}_{loc}$ , we can express the value of  $f(\Phi)$  at an arbitrary point  $\mathbf{x} \in \mathbb{R}^2$  by

$$f(\Phi)(\mathbf{x}) = \tilde{f}_{loc}((\partial_z^a \partial_{\bar{z}}^b \Phi(\mathbf{x}))_{a,b=0}^\infty). \tag{4.14}$$

Thus, the map  $f$  can be expressed, at least heuristically, by first computing various derivatives of the signal and then applying to them the invariant map  $\tilde{f}_{loc}$ , independently at each  $\mathbf{x} \in \mathbb{R}^2$ .

The expression (4.14) has the following interpretation in terms of information flow and the two different kinds of charges introduced above. Given an input signal  $\Phi \in V$  and  $\mathbf{x} \in \mathbb{R}^2$ , the signal has global charge  $\mu = 0$ , but, in general, contains multiple components having different values of the local charge  $\eta$  with respect to  $\mathbf{x}$ , according to the decomposition  $V = V_{\mu=0} = \bigoplus_{\eta \in \mathbb{Z}} V_{0,\eta}^{(\mathbf{x})}$ . By (4.11), a differential operator  $\partial_z^a \partial_{\bar{z}}^b$  maps the space  $V_{0,\eta}^{(\mathbf{x})}$  to the space  $V_{a-b,\eta+b-a}^{(\mathbf{x})}$ . However, if a signal  $\Psi \in V_{a-b,\eta+b-a}^{(\mathbf{x})}$  is continuous at  $\mathbf{x}$ , then  $\Psi$  must vanish there unless  $\eta + b - a = 0$  (see the definition (4.9), (4.10)), i.e., only information from the  $V_{0,\eta}^{(\mathbf{x})}$ -component of  $\Phi$  with  $\eta = a - b$  is observed in  $\partial_z^a \partial_{\bar{z}}^b \Phi(\mathbf{x})$ . Thus, at each point  $\mathbf{x}$ , the differential operator  $\partial_z^a \partial_{\bar{z}}^b$  can be said to transform information contained in  $\Phi$  and associated with global charge  $\mu = 0$  and local charge  $\eta = a - b$  into information associated with global charge  $\mu = a - b$  and local charge  $\eta = 0$ . This transformation is useful to us because the local charge only reflects the structure of the input signal, while the global charge is a part of the architecture of the computational model and can be used to directly control the information flow. The operators  $\partial_z^a \partial_{\bar{z}}^b$  deliver to the point  $\mathbf{x}$  information about the signal values away from this point—similarly to how this is done by local

convolutions in the convnets of Sect. 3—but now this information flow is equivariant with respect to the action of  $SO(2)$ .

By (4.14), the  $SE(2)$ -equivariant map  $f$  can be heuristically decomposed into the family of  $SE(2)$ -equivariant differentiations producing “local features”  $\partial_z^a \partial_{\bar{z}}^b \Phi(\mathbf{x})$  and followed by the  $SO(2)$ -invariant map  $\tilde{f}_{loc}$  acting independently at each  $\mathbf{x}$ . In the sequel, we use this decomposition as a general strategy in our construction of the finite convnet-like approximation model in Sect. 4.2—the “charge-conserving convnet”—and in the proof of its universality in Sect. 4.3.

The Taylor expansion (4.12) is not rigorously applicable to generic signals  $\Phi \in L^2(\mathbb{R}^2, \mathbb{R}^{d_V})$ . Therefore, we will add smoothing in our convnet-like model, to be performed before the differentiation operations. This will be discussed below in Sect. 4.1.3. Also, we will discuss there the discretization of the differential operators, in order to formulate the charge-conserving convnet as a finite computational model.

The invariant map  $\tilde{f}_{loc}$  can be approximated using invariant polynomials, as we discuss in Sect. 4.1.4 below. As discussed earlier in Sect. 2, invariant polynomials can be produced from a set of generating polynomials; however, in the present setting this set is rather large and grows rapidly as charge is increased, so it will be more efficient to just generate new invariant polynomials by multiplying general polynomials of lower degree subject to charge conservation. As a result, we will approximate the map  $\tilde{f}_{loc}$  by a series of multiplication layers in the charge-conserving convnet.

### 4.1.3 Discretized Differential Operators

Like in Sect. 3, we aim to formulate the approximation model as a computation which is fully finite except for the initial discretization of the input signal. Therefore we need to discretize the equivariant differential operators considered in Sect. 4.1.2. Given a discretized signal  $\Phi : (\lambda\mathbb{Z})^2 \rightarrow \mathbb{R}^{d_V}$  on the grid of spacing  $\lambda$ , and writing grid points as  $\gamma = (\lambda\gamma_x, \lambda\gamma_y) \in (\lambda\mathbb{Z})^2$ , we define the discrete derivatives  $\partial_z^{(\lambda)}, \partial_{\bar{z}}^{(\lambda)}$  by

$$\begin{aligned} \partial_z^{(\lambda)} \Phi(\lambda\gamma_x, \lambda\gamma_y) &= \frac{1}{4\lambda} \left( \Phi(\lambda(\gamma_x + 1), \lambda\gamma_y) - \Phi(\lambda(\gamma_x - 1), \lambda\gamma_y) \right. \\ &\quad \left. - i \left( \Phi(\lambda\gamma_x, \lambda(\gamma_y + 1)) - \Phi(\lambda\gamma_x, \lambda(\gamma_y - 1)) \right) \right), \end{aligned} \tag{4.15}$$

$$\begin{aligned} \partial_{\bar{z}}^{(\lambda)} \Phi(\lambda\gamma_x, \lambda\gamma_y) &= \frac{1}{4\lambda} \left( \Phi(\lambda(\gamma_x + 1), \lambda\gamma_y) - \Phi(\lambda(\gamma_x - 1), \lambda\gamma_y) \right. \\ &\quad \left. + i \left( \Phi(\lambda\gamma_x, \lambda(\gamma_y + 1)) - \Phi(\lambda\gamma_x, \lambda(\gamma_y - 1)) \right) \right). \end{aligned} \tag{4.16}$$

Since general signals  $\Phi \in L^2(\mathbb{R}^2, \mathbb{R}^{d_V})$  are not differentiable, we will smoothen them prior to differentiating. Smoothing will also be a part of the computational model and can be implemented by local operations as follows. Consider the discrete Laplacian

$\Delta^{(\lambda)}$  defined by

$$\begin{aligned} \Delta^{(\lambda)} \Phi(\lambda\gamma_x, \lambda\gamma_y) &= \frac{1}{\lambda^2} \left( \Phi(\lambda(\gamma_x + 1), \lambda\gamma_y) + \Phi(\lambda(\gamma_x - 1), \lambda\gamma_y) \right. \\ &\quad \left. + \Phi(\lambda\gamma_x, \lambda(\gamma_y + 1)) + \Phi(\lambda\gamma_x, \lambda(\gamma_y - 1)) - 4\Phi(\lambda\gamma_x, \lambda\gamma_y) \right). \end{aligned} \tag{4.17}$$

Then, a single smoothing layer can be implemented by the positive semidefinite operator  $1 + \frac{\lambda^2}{8} \Delta^{(\lambda)}$  :

$$\begin{aligned} \left( 1 + \frac{\lambda^2}{8} \Delta^{(\lambda)} \right) \Phi(\lambda\gamma_x, \lambda\gamma_y) &= \frac{1}{8} \left( \Phi(\lambda(\gamma_x + 1), \lambda\gamma_y) + \Phi(\lambda(\gamma_x - 1), \lambda\gamma_y) \right. \\ &\quad \left. + \Phi(\lambda\gamma_x, \lambda(\gamma_y + 1)) + \Phi(\lambda\gamma_x, \lambda(\gamma_y - 1)) \right. \\ &\quad \left. + 4\Phi(\lambda\gamma_x, \lambda\gamma_y) \right). \end{aligned} \tag{4.18}$$

We will then replace the differential operators  $\partial_z^a \partial_{\bar{z}}^b$  used in the heuristic argument in Sect. 4.1.2 by the discrete operators

$$\mathcal{L}_\lambda^{(a,b)} = (\partial_z^{(\lambda)})^a (\partial_{\bar{z}}^{(\lambda)})^b \left( 1 + \frac{\lambda^2}{8} \Delta^{(\lambda)} \right)^{\lceil 4/\lambda^2 \rceil} P_\lambda. \tag{4.19}$$

Here  $P_\lambda$  is the discretization projector (3.6). The power  $\lceil 4/\lambda^2 \rceil$  (i.e., the number of smoothing layers) scales with  $\lambda$  so that in the continuum limit  $\lambda \rightarrow 0$  the operators  $\mathcal{L}_\lambda^{(a,b)}$  converge to convolution operators. Specifically, consider the function  $\Psi_{a,b} : \mathbb{R}^2 \rightarrow \mathbb{R}$ :

$$\Psi_{a,b} = \partial_z^a \partial_{\bar{z}}^b \left( \frac{1}{2\pi} e^{-|\mathbf{x}|^2/2} \right), \tag{4.20}$$

where we identify  $|\mathbf{x}|^2 \equiv z\bar{z}$ . Define the operator  $\mathcal{L}_0^{(a,b)}$  by  $\mathcal{L}_0^{(a,b)} \Phi = \Phi * \Psi_{a,b}$ , i.e.

$$\mathcal{L}_0^{(a,b)} \Phi(\mathbf{x}) = \int_{\mathbb{R}^2} \Phi(\mathbf{x} - \mathbf{y}) \Psi_{a,b}(\mathbf{y}) d^2 \mathbf{y}. \tag{4.21}$$

Then we have the following lemma proved in Appendix A.

**Lemma 4.1** *Let  $a, b$  be fixed nonnegative integers. For all  $\lambda \in [0, 1]$ , consider the linear operators  $\mathcal{L}_\lambda^{(a,b)}$  as operators from  $L^2(\mathbb{R}^2, \mathbb{R}^{d_V})$  to  $L^\infty(\mathbb{R}^2, \mathbb{R}^{d_V})$ . Then:*

1. *The operators  $\mathcal{L}_\lambda^{(a,b)}$  are bounded uniformly in  $\lambda$ ;*
2. *As  $\lambda \rightarrow 0$ , the operators  $\mathcal{L}_\lambda^{(a,b)}$  converge strongly to the operator  $\mathcal{L}_0^{(a,b)}$ . Moreover, this convergence is uniform on compact sets  $K \subset V$  (i.e.,  $\lim_{\lambda \rightarrow 0} \sup_{\Phi \in K} \|\mathcal{L}_\lambda^{(a,b)} \Phi - \mathcal{L}_0^{(a,b)} \Phi\|_\infty = 0$ ).*

This lemma is essentially just a slight modification of Central Limit Theorem. It will be convenient to consider  $L^\infty$  rather than  $L^2$  in the target space because of the pointwise polynomial action of the layers following the smoothing and differentiation layers.

### 4.1.4 Polynomial Approximations on SO(2)-Modules

Our derivation of the approximating model in Sect. 4.1.2 was based on identifying the SO(2)-invariant map  $f_{loc}$  introduced in (4.2) and expressing it via  $\tilde{f}_{loc}$  by Eq. (4.13). It is convenient to approximate the map  $\tilde{f}_{loc}$  by invariant polynomials on appropriate SO(2)-modules, and in this section we state several general facts relevant for this purpose.

First, the following lemma is obtained immediately using symmetrization and the Weierstrass theorem (see e.g. the proof of Proposition (2.5)).

**Lemma 4.2** *Let  $f : W \rightarrow \mathbb{R}$  be a continuous SO(2)-invariant map on a real finite-dimensional SO(2)-module  $W$ . Then  $f$  can be approximated by polynomial invariants on  $W$ .*

We therefore focus on constructing general polynomial invariants on SO(2)-modules. This can be done in several ways; we will describe just one particular construction performed in a “layerwise” fashion resembling convnet layers.

It is convenient to first consider the case of SO(2)-modules over the field  $\mathbb{C}$ , since the representation theory of the group SO(2) is especially easily described when the underlying field is  $\mathbb{C}$ . Let us identify elements of SO(2) with the unit complex numbers  $e^{i\phi}$ . Then all complex irreducible representations of SO(2) are one-dimensional characters indexed by the number  $\xi \in \mathbb{Z}$ :

$$R_{e^{i\phi}\mathbf{x}} = e^{i\xi\phi}\mathbf{x}. \tag{4.22}$$

The representation  $R$  induces the dual representation acting on functions  $f(\mathbf{x})$ :

$$R_{e^{i\phi}}^* f(\mathbf{x}) = f(R_{e^{-i\phi}}\mathbf{x}).$$

In particular, if  $z_\xi$  is the variable associated with the one-dimensional space where representation (4.22) acts, then it is transformed by the dual representation as

$$R_{e^{i\phi}z_\xi}^* = e^{-i\xi\phi}z_\xi.$$

Now let  $W$  be a general finite-dimensional SO(2)-module over  $\mathbb{C}$ . Then  $W$  can be decomposed as

$$W = \bigoplus_{\xi} W_{\xi}, \tag{4.23}$$

where  $W_\xi \cong \mathbb{C}^{d_\xi}$  is the isotypic component of the representation (4.22). Let  $z_{\xi k}, k = 1, \dots, d_\xi$ , denote the variables associated with the subspace  $W_\xi$ . If  $f$  is a polynomial on  $W$ , we can write it as a linear combination of monomials:

$$f = \sum_{\mathbf{a}=(a_{\xi k})} c_{\mathbf{a}} \prod_{\xi,k} z_{\xi k}^{a_{\xi k}} \tag{4.24}$$

Then the dual representation acts on  $f$  by

$$R_{e^{i\phi}}^* f = \sum_{\mathbf{a}=(a_{\xi k})} e^{-i \sum_{\xi,k} \xi a_{\xi k} \phi} c_{\mathbf{a}} \prod_{\xi,k} z_{\xi k}^{a_{\xi k}}.$$

We see that a polynomial is invariant iff it consists of invariant monomials, and a monomial is invariant iff  $\sum_{\xi,k} \xi a_{\xi k} = 0$ .

We can generate an arbitrary  $SO(2)$ -invariant polynomial on  $W$  in the following “layer-wise” fashion. Suppose that  $\{f_{t-1,\xi,n}\}_{n=1}^{N_{t-1}}$  is a collection of polynomials generated after  $t - 1$  layers so that

$$R_{e^{i\phi}}^* f_{t-1,\xi,n} = e^{-i\xi\phi} f_{t-1,\xi,n} \tag{4.25}$$

for all  $\xi, n$ . Consider new polynomials  $\{f_{t,\xi,n}\}_{n=1}^{N_t}$  obtained from  $\{f_{t-1,\xi,n}\}_{n=1}^{N_{t-1}}$  by applying the second degree expressions

$$\begin{aligned} f_{t,\xi,n} &= w_{0,n}^{(t)} \mathbf{1}_{\xi=0} + \sum_{n_1=1}^{N_{t-1}} w_{1,\xi,n,n_1}^{(t)} f_{t-1,\xi,n_1} \\ &+ \sum_{\xi_1+\xi_2=\xi} \sum_{n_1=1}^{N_{t-1}} \sum_{n_2=1}^{N_{t-1}} w_{2,\xi_1,\xi_2,n,n_1,n_2}^{(t)} f_{t-1,\xi_1,n_1} f_{t-1,\xi_2,n_2} \end{aligned} \tag{4.26}$$

with some (complex) coefficients  $w_{0,n}^{(t)}, w_{1,\xi,n,n_1}^{(t)}, w_{2,\xi_1,\xi_2,n,n_1,n_2}^{(t)}$ . The first term is present only for  $\xi = 0$ . The third term includes the “charge conservation” constraint  $\xi = \xi_1 + \xi_2$ . It is clear that ones condition (4.25) holds for  $\{f_{t-1,\xi,n}\}_{n=1}^{N_{t-1}}$ , it also holds for  $\{f_{t,\xi,n}\}_{n=1}^{N_t}$ .

On the other hand, suppose that the initial set  $\{f_{1,\xi,n}\}_{n=1}^{N_1}$  includes all variables  $z_{\xi k}$ . Then for any invariant polynomial  $f$  on  $W$ , we can arrange the parameters  $N_t$  and the coefficients in Eq. (4.26) so that at some  $t$  we obtain  $f_{t,\xi=0,1} = f$ . Indeed, first note that thanks to the second term in Eq. (4.26) it suffices to show this for the case when  $f$  is an invariant monomial (since any invariant polynomial is a linear combination of invariant monomials, and the second term allows us to form and pass forward such linear combinations). If  $f$  is a constant, then it can be produced using the first term in Eq. (4.26). If  $f$  is a monomial of a positive degree, then it can be produced by multiplying lower degree monomials, which is afforded by the third term in Eq. (4.26).

Now we discuss the case of the underlying field  $\mathbb{R}$ . In this case, apart from the trivial one-dimensional representation, all irreducible representations of  $SO(2)$  are two-dimensional and indexed by  $\xi = 1, 2, \dots$ :

$$R_{e^{i\phi}} \begin{pmatrix} x \\ y \end{pmatrix} = \begin{pmatrix} \cos \xi \phi & \sin \xi \phi \\ -\sin \xi \phi & \cos \xi \phi \end{pmatrix} \begin{pmatrix} x \\ y \end{pmatrix}. \tag{4.27}$$

It is convenient to diagonalize such a representation, turning it into a pair of complex conjugate one-dimensional representations:

$$R_{e^{i\phi}} \begin{pmatrix} z \\ \bar{z} \end{pmatrix} = \begin{pmatrix} e^{-i\xi\phi} & 0 \\ 0 & e^{i\xi\phi} \end{pmatrix} \begin{pmatrix} z \\ \bar{z} \end{pmatrix}, \tag{4.28}$$

where

$$z = x + iy, \quad \bar{z} = x - iy.$$

More generally, any real  $SO(2)$ -module  $W$  can be decomposed exactly as in (4.23) into isotypic components  $W_\xi$  associated with complex characters, but with the additional constraints

$$W_\xi = \overline{W_{-\xi}}, \tag{4.29}$$

meaning that  $d_\xi = d_{-\xi}$  and

$$W_\xi = W_{\xi, \text{Re}} + iW_{\xi, \text{Im}}, \quad W_{-\xi} = W_{\xi, \text{Re}} - iW_{\xi, \text{Im}}, \quad (\xi = 1, 2, \dots)$$

with some real  $d_\xi$ -dimensional spaces  $W_{\xi, \text{Re}}, W_{\xi, \text{Im}}$ .

Any polynomial on  $W$  can then be written in terms of real variables  $z_{0,k}$  corresponding to  $\xi = 0$  and complex variables

$$z_{\xi,k} = x_{\xi k} + iy_{\xi k}, \quad z_{-\xi,k} = x_{\xi k} - iy_{\xi k} \quad (\xi = 1, 2, \dots) \tag{4.30}$$

constrained by the relations

$$z_{\xi,k} = \overline{z_{-\xi,k}}.$$

Suppose that a polynomial  $f$  on  $W$  is expanded over monomials in  $z_{\xi,k}$  as in Eq. (4.24). This expansion is unique (the coefficients are given by

$$c_{\mathbf{a}} = \left( \prod_{\xi,k} \frac{\partial^{a_{\xi,k}}}{a_{\xi,k}!} \right) f(0),$$

where  $\partial_{z_{\xi,k}} = \frac{1}{2}(\partial_{x_{\xi k}} - i\partial_{y_{\xi k}})$  for  $\xi > 0$  and  $\partial_{z_{\xi,k}} = \frac{1}{2}(\partial_{x_{-\xi,k}} + i\partial_{y_{-\xi,k}})$  for  $\xi < 0$ ). This implies that the condition for the polynomial  $f$  to be invariant on  $W$  is the same

as in the previously considered complex case: the polynomial must consist of invariant monomials, and a monomial is invariant iff  $\sum_{\xi,k} \xi a_{\xi k} = 0$ .

Therefore, in the case of real  $SO(2)$ -modules, any invariant polynomial can be generated using the same procedure described earlier for the complex case, i.e., by taking the complex extension of the module and iteratively generating (complex) polynomials  $\{f_{t,\xi,n}\}_{n=1}^{N_t}$  using Eq. (4.26). The real part of a complex invariant polynomial on a real module is a real invariant polynomial. Thus, to ensure that in the case of real modules  $W$  the procedure produces all real invariant polynomials, and only such polynomials, we can just add taking the real part of  $f_{t,\xi=0,1}$  at the last step of the procedure.

### 4.2 Charge-Conserving Convnet

We can now describe precisely our convnet-like model for approximating arbitrary  $SE(2)$ -equivariant continuous maps  $f : V \rightarrow U$ , where  $V = L^2(\mathbb{R}^2, \mathbb{R}^{d_V})$ ,  $U = L^2(\mathbb{R}^2, \mathbb{R}^{d_U})$ . The overview of the model is given in Fig. 3. Like the models of Sect. 3, the present model starts with the discretization projection followed by some finite computation. The model includes three groups of layers: smoothing layers ( $\mathcal{L}_{\text{smooth}}$ ), differentiation layers ( $\mathcal{L}_{\text{diff}}$ ) and multiplication layers ( $\mathcal{L}_{\text{mult}}$ ). The parameters of the model are the lattice spacing  $\lambda$ , cutoff range  $\Lambda$  of the output, dimension  $d_{\text{mult}}$  of auxiliary spaces, and the numbers  $T_{\text{diff}}, T_{\text{mult}}$  of differentiation and multiplication layers. The overall operation of the model can be described as the chain

$$\widehat{f} : V \xrightarrow{P_{\lambda,\Lambda'}} V_{\lambda,\Lambda'} (\equiv W_1) \xrightarrow{\mathcal{L}_{\text{smooth}}} W_{\text{smooth}} \xrightarrow{\mathcal{L}_{\text{diff}}} W_{\text{diff}} \xrightarrow{\mathcal{L}_{\text{mult}}} U_{\lambda,\Lambda}. \tag{4.31}$$

We describe now all these layers in detail.

**Initial projection** The initial discretization projection  $P_{\lambda,\Lambda'}$  is defined as explained in Sect. 3 after Eq. (3.8). The input cutoff range  $\Lambda'$  is given by  $\Lambda' = \Lambda + (T_{\text{diff}} + \lceil 4/\lambda^2 \rceil)\lambda$ . This padding ensures that the output cutoff range will be equal to the specified value  $\Lambda$ . With respect to the spatial grid structure, the space  $W_1$  can be decomposed as

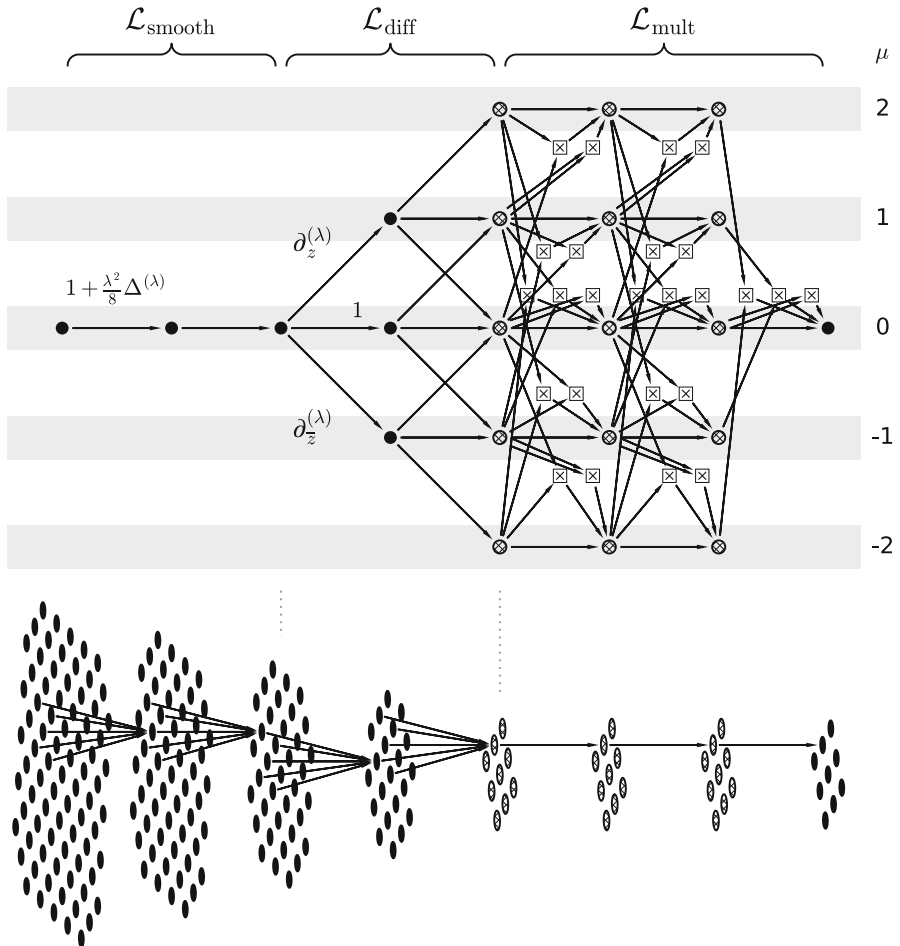
$$W_1 = \bigoplus_{\gamma \in \lambda Z_{\lfloor \Lambda'/\lambda \rfloor}} \mathbb{R}^{d_V},$$

where  $Z_L$  is the cubic subset of the grid defined in (3.7).

**Smoothing layers** The model contains  $\lceil 4/\lambda^2 \rceil$  smoothing layers performing the same elementary smoothing operation  $1 + \frac{\lambda^2}{8} \Delta^{(\lambda)}$ :

$$\mathcal{L}_{\text{smooth}} = \left( 1 + \frac{\lambda^2}{8} \Delta^{(\lambda)} \right)^{\lceil 4/\lambda^2 \rceil},$$

where the discrete Laplacian  $\Delta^{(\lambda)}$  is defined as in Eq. (4.17). In each layer the value of the transformed signal at the current spatial position is determined by the values of the signal in the previous layer at this position and its 4 nearest neighbors as given



**Fig. 3** Architecture of the charge-conserving convnet. The top figure shows the information flow in the fixed-charge subspaces of the feature space, while the bottom figure shows the same flow in the spatial coordinates. The smoothing layers only act on spatial dimensions, the multiplication layers only on feature dimensions, and the differentiation layers both on spatial and feature dimensions. Operation of smoothing and differentiation layers only involves nearest neighbors while in the multiplication layers the transitions are constrained by the requirement of charge conservation. The smoothing and differentiation layers are linear; the multiplication layers are not. The last multiplication layer only has zero-charge (SO(2)-invariant) output

in Eq. (4.18). Accordingly, the domain size shrinks with each layer so that the output space of  $\mathcal{L}_{\text{smooth}}$  can be written as



$$W_{\text{smooth}} = \bigoplus_{\gamma \in \lambda Z_{\lfloor \Lambda''/\lambda \rfloor}} \mathbb{R}^{d_V},$$

where  $\Lambda'' = \Lambda' - \lceil 4/\lambda^2 \rceil \lambda = \Lambda + T_{\text{diff}} \lambda$ .

**Differentiation layers** The model contains  $T_{\text{diff}}$  differentiation layers computing the discretized derivatives  $\partial_z^{(\lambda)}$ ,  $\partial_{\bar{z}}^{(\lambda)}$  as defined in (4.15), (4.16). Like the smoothing layers, these derivatives shrink the domain, but additionally, as discussed in Sect. 4.1.2, they change the representation of the group  $\text{SE}(2)$  associated with the global charge  $\mu$  (see Eq. (4.7)).

Denoting the individual differentiation layers by  $\mathcal{L}_{\text{diff},t}$ ,  $t = 1, \dots, T_{\text{diff}}$ , their action can be described as the chain

$$\mathcal{L}_{\text{diff}} : W_{\text{smooth}} \xrightarrow{\mathcal{L}_{\text{diff},1}} W_{\text{diff},1} \xrightarrow{\mathcal{L}_{\text{diff},2}} W_{\text{diff},2} \dots \xrightarrow{\mathcal{L}_{\text{diff},T_{\text{diff}}}} W_{\text{diff},T_{\text{diff}}} (\equiv W_{\text{diff}}).$$

We decompose each intermediate space  $W_{\text{diff},t}$  into subspaces characterized by degree  $s$  of the derivative and by charge  $\mu$ :

$$W_{\text{diff},t} = \bigoplus_{s=0}^t \bigoplus_{\mu=-s}^s W_{\text{diff},t,s,\mu}. \tag{4.32}$$

Each  $W_{\text{diff},t,s,\mu}$  can be further decomposed as a direct sum over the grid points:

$$W_{\text{diff},t,s,\mu} = \bigoplus_{\gamma \in \lambda Z_{\lfloor \Lambda/\lambda \rfloor + T_{\text{diff}} - t}} \mathbb{C}^{d_V}. \tag{4.33}$$

Consider the operator  $L_{\text{diff},t}$  as a block matrix with respect to decomposition (4.32) of the input and output spaces  $W_{\text{diff},t-1}$ ,  $W_{\text{diff},t}$ , and denote by  $(\mathcal{L}_{\text{diff},t})_{(s_{t-1}, \mu_{t-1}) \rightarrow (s_t, \mu_t)}$  the respective blocks. Then we define

$$(\mathcal{L}_{\text{diff},t})_{(s_{t-1}, \mu_{t-1}) \rightarrow (s_t, \mu_t)} = \begin{cases} \partial_z^{(\lambda)}, & \text{if } s_t = s_{t-1} + 1, \mu_t = \mu_{t-1} + 1, \\ \partial_{\bar{z}}^{(\lambda)}, & \text{if } s_t = s_{t-1} + 1, \mu_t = \mu_{t-1} - 1, \\ \mathbf{1}, & \text{if } s_t = s_{t-1}, \mu_t = \mu_{t-1}, \\ 0, & \text{otherwise.} \end{cases} \tag{4.34}$$

With this definition, the final space  $W_{\text{diff},T_{\text{diff}}}$  contains all discrete derivatives  $(\partial_z^{(\lambda)})^a (\partial_{\bar{z}}^{(\lambda)})^b \Phi$  of the smoothed signal  $\Phi \in W_{\text{smooth}}$  of degrees  $s = a + b \leq T_{\text{diff}}$ . Each such derivative can be obtained by arranging the elementary steps (4.34) in different order, so that the derivative will actually appear in  $W_{\text{diff},T_{\text{diff}}}$  with the coefficient  $\frac{T_{\text{diff}}!}{a!b!(T_{\text{diff}}-a-b)!}$ . This coefficient is not important for the subsequent exposition.

**Multiplication layers** In contrast to the smoothing and differentiation layers, the multiplication layers act strictly locally (pointwise). These layers implement products and linear combinations of signals of the preceding layers subject to conservation of global charge, based on the procedure of generation of invariant polynomials described in Sect. 4.1.4.

Denoting the individual layers by  $\mathcal{L}_{\text{mult},t}$ ,  $t = 1, \dots, T_{\text{mult}}$ , their action is described by the chain

$$\mathcal{L}_{\text{mult}} : W_{\text{diff}} \xrightarrow{\mathcal{L}_{\text{mult},1}} W_{\text{mult},1} \xrightarrow{\mathcal{L}_{\text{mult},2}} W_{\text{mult},2} \dots \xrightarrow{\mathcal{L}_{\text{mult},T_{\text{mult}}}} W_{\text{mult},T_{\text{mult}}} \equiv U_{\lambda,\Lambda}.$$

Each space  $W_{\text{mult},t}$  except for the final one ( $W_{\text{mult},T_{\text{mult}}}$ ) is decomposed into subspaces characterized by spatial position  $\gamma \in (\lambda\mathbb{Z})^2$  and charge  $\mu$ :

$$W_{\text{mult},t} = \bigoplus_{\gamma \in \lambda Z_{[\Lambda/\lambda]}} \bigoplus_{\mu=-T_{\text{diff}}}^{T_{\text{diff}}} W_{\text{mult},t,\gamma,\mu}. \tag{4.35}$$

Each space  $W_{\text{mult},t,\gamma,\mu}$  is a complex  $d_{\text{mult}}$ -dimensional space, where  $d_{\text{mult}}$  is a parameter of the model:

$$W_{\text{mult},t,\gamma,\mu} = \mathbb{C}^{d_{\text{mult}}}.$$

The final space  $W_{\text{mult},T_{\text{mult}}}$  is real,  $d_U$ -dimensional, and only has the charge-0 component:

$$W_{\text{mult},T_{\text{mult}}} = \bigoplus_{\gamma \in \lambda Z_{[\Lambda/\lambda]}} W_{\text{mult},T_{\text{mult}},\gamma,\mu=0}, \quad W_{\text{mult},T_{\text{mult}},\gamma,\mu=0} = \mathbb{R}^{d_U},$$

so that  $W_{\text{mult},T_{\text{mult}}}$  can be identified with  $U_{\lambda,\Lambda}$ . The initial space  $W_{\text{diff}}$  can also be expanded in the form (4.35) by reshaping its components (4.32), (4.33):

$$\begin{aligned} W_{\text{diff}} &= \bigoplus_{s=0}^{T_{\text{diff}}} \bigoplus_{\mu=-s}^s W_{\text{diff},T_{\text{diff}},s,\mu} \\ &= \bigoplus_{s=0}^{T_{\text{diff}}} \bigoplus_{\mu=-s}^s \bigoplus_{\gamma \in \lambda Z_{[\Lambda/\lambda]}} \mathbb{C}^{d_V} \\ &= \bigoplus_{\gamma \in \lambda Z_{[\Lambda/\lambda]}} \bigoplus_{\mu=-T_{\text{diff}}}^{T_{\text{diff}}} W_{\text{mult},0,\gamma,\mu}, \end{aligned}$$

where

$$W_{\text{mult},0,\gamma,\mu} = \bigoplus_{s=|\mu|}^{T_{\text{diff}}} \mathbb{C}^{d_V}.$$

The multiplication layers  $\mathcal{L}_{\text{mult},t}$  act separately and identically at each  $\gamma \in \lambda Z_{[\Lambda/\lambda]}$ , i.e., without loss of generality these layers can be thought of as maps

$$\mathcal{L}_{\text{mult},t} : \bigoplus_{\mu=-T_{\text{diff}}}^{T_{\text{diff}}} W_{\text{mult},t-1,\gamma=0,\mu} \longrightarrow \bigoplus_{\mu=-T_{\text{diff}}}^{T_{\text{diff}}} W_{\text{mult},t,\gamma=0,\mu}.$$

To define  $\mathcal{L}_{\text{mult},t}$ , let us represent its input  $\Phi \in \bigoplus_{\mu=-T_{\text{diff}}}^{T_{\text{diff}}} W_{\text{mult},t-1,\gamma=0,\mu}$  as

$$\Phi = \sum_{\mu=-T_{\text{diff}}}^{T_{\text{diff}}} \Phi_{\mu} = \sum_{\mu=-T_{\text{diff}}}^{T_{\text{diff}}} \sum_{n=1}^{d_{\text{mult}}} \Phi_{\mu,n} \mathbf{e}_{\mu,n},$$

where  $\mathbf{e}_{\mu,n}$  denote the basis vectors in  $W_{\text{mult},t-1,\gamma=0,\mu}$ . We represent the output  $\Psi \in \bigoplus_{\mu=-T_{\text{diff}}}^{T_{\text{diff}}} W_{\text{mult},t,\gamma=0,\mu}$  of  $\mathcal{L}_{\text{mult},t}$  in the same way:

$$\Psi = \sum_{\mu=-T_{\text{diff}}}^{T_{\text{diff}}} \Psi_{\mu} = \sum_{\mu=-T_{\text{diff}}}^{T_{\text{diff}}} \sum_{n=1}^{d_{\text{mult}}} \Psi_{\mu,n} \mathbf{e}_{\mu,n}.$$

Then, based on Eq. (4.26), for  $t < T_{\text{mult}}$  we define  $\mathcal{L}_{\text{mult},t} \Phi = \Psi$  by

$$\begin{aligned} \Psi_{\mu,n} &= w_{0,n}^{(t)} \mathbf{1}_{\mu=0} + \sum_{n_1=1}^{d_{\text{mult}}} w_{1,\mu,n,n_1}^{(t)} \Phi_{\mu,n_1} \\ &+ \sum_{\substack{-T_{\text{diff}} \leq \mu_1, \mu_2 \leq T_{\text{diff}} \\ \mu_1 + \mu_2 = \mu}} \sum_{n_1=1}^{d_{\text{mult}}} \sum_{n_2=1}^{d_{\text{mult}}} w_{2,\mu_1,\mu_2,n,n_1,n_2}^{(t)} \Phi_{\mu_1,n_1} \Phi_{\mu_2,n_2}, \end{aligned} \quad (4.36)$$

with some complex weights  $w_{0,n}^{(t)}, w_{1,\mu,n,n_1}^{(t)}, w_{2,\mu_1,\mu_2,n,n_1,n_2}^{(t)}$ . In the final layer  $t = T_{\text{mult}}$  the network only needs to generate a real charge-0 (invariant) vector, so in this case  $\Psi$  only has real  $\mu = 0$  components:

$$\begin{aligned} \Psi_{0,n} &= \text{Re} \left( w_{0,n}^{(t)} + \sum_{n_1=1}^{d_{\text{mult}}} w_{1,0,n,n_1}^{(t)} \Phi_{0,n_1} \right. \\ &\left. + \sum_{\substack{-T_{\text{diff}} \leq \mu_1, \mu_2 \leq T_{\text{diff}} \\ \mu_1 + \mu_2 = 0}} \sum_{n_1=1}^{d_{\text{mult}}} \sum_{n_2=1}^{d_{\text{mult}}} w_{2,\mu_1,\mu_2,n,n_1,n_2}^{(t)} \Phi_{\mu_1,n_1} \Phi_{\mu_2,n_2} \right). \end{aligned} \quad (4.37)$$

This completes the description of the charge-conserving convnet. In the sequel, it will be convenient to consider a family of convnets having all parameters and weights in common except for the grid spacing  $\lambda$ . Observe that this parameter can be varied independently of all other parameters and weights  $(\Lambda, d_{\text{mult}}, T_{\text{diff}}, T_{\text{mult}}, w_{0,n}^{(t)}, w_{1,\mu,n,n_1}^{(t)}, w_{2,\mu_1,\mu_2,n,n_1,n_2}^{(t)})$ . The parameter  $\lambda$  affects the number of smoothing layers, and decreasing this parameter means that essentially the same convnet is applied at a higher resolution. Accordingly, we will call such a family a “multi-resolution convnet”.

**Definition 4.1** A **charge-conserving convnet** is a map  $\widehat{f} : V \rightarrow U$  given in (4.31), characterized by parameters  $\lambda, \Lambda, d_{\text{mult}}, T_{\text{diff}}, T_{\text{mult}}$  and weights  $w_{0,n}^{(t)}, w_{1,\mu,n,n_1}^{(t)}, w_{2,\mu_1,\mu_2,n,n_1,n_2}^{(t)}$ , and constructed as described above. A **multi-resolution charge-conserving convnet**  $\widehat{f}_{\lambda}$  is obtained by arbitrarily varying the grid spacing parameter  $\lambda$  in the charge-conserving convnet  $\widehat{f}$ .

We comment now why it is natural to call this model “charge-conserving”. As already explained in Sect. 4.1.2, if the intermediate spaces labeled by specific  $\mu$ ’s are equipped with the special representations (4.6), then, up to the spatial cutoff, the

differentiation layers  $\mathcal{L}_{\text{diff}}$  are SE(2)-equivariant and conserve the “total charge”  $\mu + \eta$ , where  $\eta$  is the “local charge” (see Eq. (4.11)). Clearly, the same can be said about the smoothing layers  $\mathcal{L}_{\text{smooth}}$  which, in fact, separately conserve the global charge  $\mu$  and the local charge  $\eta$ . Moreover, observe that the multiplication layers  $\mathcal{L}_{\text{mult}}$ , though nonlinear, are also equivariant and separately conserve the charges  $\mu$  and  $\eta$ . Indeed, consider the transformations (4.36), (4.37). The first term in these transformations creates an SE(2)-invariant,  $\mu = \eta = 0$  signal. The second, linear term does not change  $\mu$  or  $\eta$  of the input signal. The third term creates products  $\Psi_\mu = \Phi_{\mu_1} \Phi_{\mu_2}$ , where  $\mu = \mu_1 + \mu_2$ . This multiplication operation is equivariant with respect to the respective representations  $R^{(\mu)}$ ,  $R^{(\mu_1)}$ ,  $R^{(\mu_2)}$  as defined in (4.6). Also, if the signals  $\Phi_{\mu_1}$ ,  $\Phi_{\mu_2}$  have local charges  $\eta_1, \eta_2$  at a particular point  $\mathbf{x}$ , then the product  $\Phi_{\mu_1} \Phi_{\mu_2}$  has local charge  $\eta = \eta_1 + \eta_2$  at this point (see Eqs.(4.9), (4.10)).

### 4.3 The Main Result

To state our main result, we define a limit point of charge-conserving convnets.

**Definition 4.2** With  $V = L^2(\mathbb{R}^2, \mathbb{R}^{d_V})$  and  $U = L^2(\mathbb{R}^2, \mathbb{R}^{d_U})$ , we say that a map  $f : V \rightarrow U$  is a **limit point of charge-conserving convnets** if for any compact set  $K \subset V$ , any  $\epsilon > 0$  and  $\Lambda_0 > 0$  there exist a multi-resolution charge-conserving convnet  $\widehat{f}_\lambda$  with  $\Lambda > \Lambda_0$  such that  $\sup_{\Phi \in K} \|\widehat{f}_\lambda(\Phi) - f(\Phi)\| \leq \epsilon$  for all sufficiently small grid spacings  $\lambda$ .

Then our main result is the following theorem.

**Theorem 4.1** Let  $V = L^2(\mathbb{R}^2, \mathbb{R}^{d_V})$  and  $U = L^2(\mathbb{R}^2, \mathbb{R}^{d_U})$ . A map  $f : V \rightarrow U$  is a limit point of charge-conserving convnets if and only if  $f$  is SE(2)-equivariant and continuous in the norm topology.

**Proof** To simplify the exposition, we will assume that  $d_V = d_U = 1$ ; generalization of all the arguments to vector-valued input and output signals is straightforward.

We start by observing that a multi-resolution family of charge-conserving convnets has a natural scaling limit as the lattice spacing  $\lambda \rightarrow 0$ :

$$\widehat{f}_0(\Phi) = \lim_{\lambda \rightarrow 0} \widehat{f}_\lambda(\Phi). \tag{4.38}$$

Indeed, by (4.31), at  $\lambda > 0$  we can represent the convnet as the composition of maps

$$\widehat{f}_\lambda = \mathcal{L}_{\text{mult}} \circ \mathcal{L}_{\text{diff}} \circ \mathcal{L}_{\text{smooth}} \circ P_{\lambda, \Lambda'}.$$

The part  $\mathcal{L}_{\text{diff}} \circ \mathcal{L}_{\text{smooth}} \circ P_{\lambda, \Lambda'}$  of this computation implements several maps  $\mathcal{L}_\lambda^{(a,b)}$  introduced in (4.19). More precisely, by the definition of differentiation layers in Sect. 4.2, the output space  $W_{\text{diff}}$  of the linear operator  $\mathcal{L}_{\text{diff}} \circ \mathcal{L}_{\text{smooth}} \circ P_{\lambda, \Lambda'}$  can be decomposed into the direct sum (4.32) over several degrees  $s$  and charges  $\mu$ . The respective components of  $\mathcal{L}_{\text{diff}} \circ \mathcal{L}_{\text{smooth}} \circ P_{\lambda, \Lambda'}$  are, up to unimportant combinatoric

coefficients, just the operators  $\mathcal{L}_\lambda^{(a,b)}$  with  $a + b = s, a - b = \mu$ :

$$\mathcal{L}_{\text{diff}} \circ \mathcal{L}_{\text{smooth}} \circ P_{\lambda, \Lambda'} = (\dots, c_{a,b} \mathcal{L}_\lambda^{(a,b)}, \dots), \quad c_{a,b} = \frac{T_{\text{diff}}!}{a!b!(T_{\text{diff}} - a - b)!}, \quad (4.39)$$

with the caveat that the output of  $\mathcal{L}_\lambda^{(a,b)}$  is spatially restricted to the bounded domain  $[-\Lambda, \Lambda]^2$ . By Lemma 4.1, as  $\lambda \rightarrow 0$ , the operators  $\mathcal{L}_\lambda^{(a,b)}$  converge to the operator  $\mathcal{L}_0^{(a,b)}$  defined in Eq. (4.21), so that for any  $\Phi \in L^2(\mathbb{R}^2)$  the signals  $\mathcal{L}_\lambda^{(a,b)} \Phi$  are bounded functions on  $\mathbb{R}^2$  and converge to  $\mathcal{L}_0^{(a,b)} \Phi$  in the uniform norm  $\|\cdot\|_\infty$ . Let us denote the limiting linear operator by  $\mathcal{L}_{\text{conv}}$ :

$$\mathcal{L}_{\text{conv}} = \lim_{\lambda \rightarrow 0} \mathcal{L}_{\text{diff}} \circ \mathcal{L}_{\text{smooth}} \circ P_{\lambda, \Lambda'}. \quad (4.40)$$

The full limiting map  $\widehat{f}_0(\Phi)$  is then obtained by pointwise application (separately at each point  $\mathbf{x} \in [-\Lambda, \Lambda]^2$ ) of the multiplication layers  $\mathcal{L}_{\text{mult}}$  to the signals  $\mathcal{L}_{T_{\text{diff}}} \Phi$ :

$$\widehat{f}_0(\Phi) = \mathcal{L}_{\text{mult}}(\mathcal{L}_{\text{conv}} \Phi). \quad (4.41)$$

For any  $\Phi \in L^2(\mathbb{R}^2)$ , this  $f_0(\Phi)$  is a well-defined bounded signal on the domain  $[-\Lambda, \Lambda]^2$ . It is bounded because the multiplication layers  $\mathcal{L}_{\text{mult}}$  implement a continuous (polynomial) map, and because, as already mentioned,  $\mathcal{L}_{\text{conv}} \Phi$  is a bounded signal. Since the domain  $[-\Lambda, \Lambda]^2$  has a finite Lebesgue measure, we have  $f_0(\Phi) \in L^\infty([-\Lambda, \Lambda]^2) \subset L^2([-\Lambda, \Lambda]^2)$ . By a similar argument, the convergence in (4.38) can be understood in the  $L^\infty([-\Lambda, \Lambda]^2)$  or  $L^2([-\Lambda, \Lambda]^2)$  sense, e.g.:

$$\|\widehat{f}_0(\Phi) - \widehat{f}_\lambda(\Phi)\|_{L^2([-\Lambda, \Lambda]^2)} \xrightarrow{\lambda \rightarrow 0} 0, \quad \Phi \in V. \quad (4.42)$$

Below, we will use the scaling limit  $\widehat{f}_0$  as an intermediate approximator.

We will now prove the necessity and then the sufficiency parts of the theorem.

**Necessity** (*a limit point  $f$  is continuous and SE(2)-equivariant*). As in the previous theorems 3.1, 3.2, continuity of  $f$  follows by standard topological arguments, and we only need to prove the SE(2)-equivariance.

Let us first prove the  $\mathbb{R}^2$ -equivariance of  $f$ . By the definition of a limit point, for any  $\Phi \in V, \mathbf{x} \in \mathbb{R}^2, \epsilon > 0$  and  $\Lambda_0 > 0$  there is a multi-resolution convnet  $\widehat{f}_\lambda$  with  $\Lambda > \Lambda_0$  such that

$$\|\widehat{f}_\lambda(\Phi) - f(\Phi)\| \leq \epsilon, \quad \|\widehat{f}_\lambda(R_{(\mathbf{x},1)} \Phi) - f(R_{(\mathbf{x},1)} \Phi)\| \leq \epsilon \quad (4.43)$$

for all sufficiently small  $\lambda$ . Consider the scaling limit  $\widehat{f}_0 = \lim_{\lambda \rightarrow 0} \widehat{f}_\lambda$  constructed above. As shown above,  $\widehat{f}_\lambda(\Phi)$  converges to  $\widehat{f}_0(\Phi)$  in the  $L^2$  sense, so the inequalities (4.43) remain valid for  $\widehat{f}_0(\Phi)$ :

$$\|\widehat{f}_0(\Phi) - f(\Phi)\| \leq \epsilon, \quad \|\widehat{f}_0(R_{(\mathbf{x},1)} \Phi) - f(R_{(\mathbf{x},1)} \Phi)\| \leq \epsilon. \quad (4.44)$$

The map  $\widehat{f}_0$  is not  $\mathbb{R}^2$ -equivariant only because its output is restricted to the domain  $[-\Lambda, \Lambda]^2$ , since otherwise both maps  $\mathcal{L}_{\text{mult}}, \mathcal{L}_{\text{conv}}$  appearing in the superposition (4.41) are  $\mathbb{R}^2$ -equivariant. Therefore, for any  $\mathbf{y} \in \mathbb{R}^2$ ,

$$\widehat{f}_0(R_{(\mathbf{x},1)}\Phi)(\mathbf{y}) = R_{(\mathbf{x},1)}\widehat{f}_0(\Phi)(\mathbf{y}) = \widehat{f}_0(\Phi)(\mathbf{y} - \mathbf{x}), \quad \text{if } \mathbf{y}, \mathbf{y} - \mathbf{x} \in [-\Lambda, \Lambda]^2. \tag{4.45}$$

Consider the set

$$\Pi_{\Lambda, \mathbf{x}} = \{\mathbf{y} \in \mathbb{R}^2 : \mathbf{y}, \mathbf{y} - \mathbf{x} \in [-\Lambda, \Lambda]^2\} = [-\Lambda, \Lambda]^2 \cap \mathcal{A}_{(\mathbf{x},1)}([-\Lambda, \Lambda]^2).$$

The identity (4.45) implies that

$$P_{\Pi_{\Lambda, \mathbf{x}}}\widehat{f}_0(R_{(\mathbf{x},1)}\Phi) = P_{\Pi_{\Lambda, \mathbf{x}}}R_{(\mathbf{x},1)}\widehat{f}_0(\Phi), \tag{4.46}$$

where  $P_{\Pi_{\Lambda, \mathbf{x}}}$  denotes the projection to the subspace  $L^2(\Pi_{\Lambda, \mathbf{x}})$  in  $L^2(\mathbb{R}^2)$ . For a fixed  $\mathbf{x}$ , the projectors  $P_{\Pi_{\Lambda, \mathbf{x}}}$  converge strongly to the identity as  $\Lambda \rightarrow \infty$ , therefore we can choose  $\Lambda$  sufficiently large so that

$$\|P_{\Pi_{\Lambda, \mathbf{x}}}f(\Phi) - f(\Phi)\| \leq \epsilon, \quad \|P_{\Pi_{\Lambda, \mathbf{x}}}f(R_{(\mathbf{x},1)}\Phi) - f(R_{(\mathbf{x},1)}\Phi)\| \leq \epsilon. \tag{4.47}$$

Then, assuming that the approximating convnet has a sufficiently large range  $\Lambda$ , we have

$$\begin{aligned} \|f(R_{(\mathbf{x},1)}\Phi) - R_{(\mathbf{x},1)}f(\Phi)\| &\leq \|f(R_{(\mathbf{x},1)}\Phi) - P_{\Pi_{\Lambda, \mathbf{x}}}f(R_{(\mathbf{x},1)}\Phi)\| \\ &\quad + \|P_{\Pi_{\Lambda, \mathbf{x}}}f(R_{(\mathbf{x},1)}\Phi) - P_{\Pi_{\Lambda, \mathbf{x}}}\widehat{f}_0(R_{(\mathbf{x},1)}\Phi)\| \\ &\quad + \|P_{\Pi_{\Lambda, \mathbf{x}}}\widehat{f}_0(R_{(\mathbf{x},1)}\Phi) - P_{\Pi_{\Lambda, \mathbf{x}}}R_{(\mathbf{x},1)}\widehat{f}_0(\Phi)\| \\ &\quad + \|P_{\Pi_{\Lambda, \mathbf{x}}}R_{(\mathbf{x},1)}\widehat{f}_0(\Phi) - P_{\Pi_{\Lambda, \mathbf{x}}}R_{(\mathbf{x},1)}f(\Phi)\| \\ &\quad + \|P_{\Pi_{\Lambda, \mathbf{x}}}R_{(\mathbf{x},1)}f(\Phi) - R_{(\mathbf{x},1)}f(\Phi)\| \\ &\leq 4\epsilon, \end{aligned}$$

where we used the bounds (4.44), (4.47), the equalities  $\|P_{\Pi_{\Lambda, \mathbf{x}}}\| = \|R_{(\mathbf{x},1)}\| = 1$ , and the identity (4.46). Taking the limit  $\epsilon \rightarrow 0$ , we obtain the desired  $\mathbb{R}^2$ -equivariance of  $f$ .

To complete the proof of SE(2)-equivariance, we will show that for any  $\theta \in \text{SO}(2)$  we have

$$R_{(0,\theta)}\widehat{f}_0(\Phi)(\mathbf{x}) = \widehat{f}_0(R_{(0,\theta)}\Phi)(\mathbf{x}), \quad \mathbf{x} \in \Pi_{\Lambda, \theta}, \tag{4.48}$$

where

$$\Pi_{\Lambda, \theta} = [-\Lambda, \Lambda]^2 \cap \mathcal{A}_{(0,\theta)}([-\Lambda, \Lambda]^2).$$

Identity (4.48) is an analog of identity (4.45) that we used to prove the  $R_{(\mathbf{x},1)}$ -equivariance of  $f$ . Once Eq. (4.48) is established, we can prove the  $R_{(0,\theta)}$ -equivariance

of  $f$  by arguing in the same way as we did above to prove the  $R_{(\mathbf{x},1)}$ -equivariance. After that, the  $R_{(0,\theta)}$ -equivariance and the  $R_{(\mathbf{x},1)}$ -equivariance together imply the full SE(2)-equivariance.

Note that by using the partial translation equivariance (4.45) and repeating the computation from Lemma 4.2, it suffices to prove the identity (4.48) only in the special case  $\mathbf{x} = 0$ :

$$\widehat{f}_0(R_{(0,\theta)}\Phi)(0) = \widehat{f}_0(\Phi)(0). \tag{4.49}$$

Indeed, suppose that Eq. (4.49) is established and  $\Lambda$  is sufficiently large so that  $\mathbf{x}, \theta^{-1}\mathbf{x} \in [-\Lambda, \Lambda]^2$ . Then,

$$\begin{aligned} \widehat{f}_0(R_{(0,\theta)}\Phi)(\mathbf{x}) &= R_{(-\mathbf{x},1)}\widehat{f}_0(R_{(0,\theta)}\Phi)(0) \\ &= \widehat{f}_0(R_{(-\mathbf{x},1)}R_{(0,\theta)}\Phi)(0) \\ &= \widehat{f}_0(R_{(0,\theta)}R_{(-\theta^{-1}\mathbf{x},1)}\Phi)(0) \\ &= \widehat{f}_0(R_{(-\theta^{-1}\mathbf{x},1)}\Phi)(0) \\ &= R_{(-\theta^{-1}\mathbf{x},1)}\widehat{f}_0(\Phi)(0) \\ &= R_{(\mathbf{x},\theta)}R_{(-\theta^{-1}\mathbf{x},1)}\widehat{f}_0(\Phi)(\mathcal{A}_{(\mathbf{x},\theta)}0) \\ &= R_{(0,\theta)}\widehat{f}_0(\Phi)(\mathbf{x}), \end{aligned}$$

where we used general properties of the representaton  $R$  (steps 1, 6, 7), Eq. (4.49) (step 4), and the partial  $\mathbb{R}^2$ -equivariance (4.45) (steps 2 and 5, using the fact that  $0, \mathbf{x}, \theta^{-1}\mathbf{x} \in [-\Lambda, \Lambda]^2$ ).

To establish Eq. (4.49), recall that, by Eq. (4.41), the value  $\widehat{f}_0(\Phi)(0)$  is obtained by first evaluating  $\mathcal{L}_{\text{conv}}(\Phi)$  at  $\mathbf{x} = 0$  and then applying to the resulting values the map  $\mathcal{L}_{\text{mult}}$ . By Eqs.(4.39), (4.40) and Lemma 4.1, we can write  $\mathcal{L}_{\text{conv}}(\Phi)(0)$  as a vector with components

$$\mathcal{L}_{\text{conv}}(\Phi)(0) = (\dots, c_{a,b}\mathcal{L}_0^{(a,b)}(\Phi)(0), \dots), \tag{4.50}$$

where, by Eq. (4.21),

$$\mathcal{L}_0^{(a,b)}\Phi(0) = \int_{\mathbb{R}^2} \Phi(-\mathbf{y})\Psi_{a,b}(\mathbf{y})d^2\mathbf{y},$$

and  $\Psi_{a,b}$  is given by Eq. (4.20):

$$\Psi_{a,b} = \partial_z^a \partial_{\bar{z}}^b \left( \frac{1}{2\pi} e^{-|\mathbf{x}|^2/2} \right) = \partial_z^a \partial_{\bar{z}}^b \left( \frac{1}{2\pi} e^{-z\bar{z}/2} \right).$$

In the language of Sect. 4.1.2,  $\Psi_{a,b}$  has local charge  $\eta = b - a$ :

$$\Psi_{a,b}(\mathcal{A}_{(0,e^{-i\phi})}\mathbf{x}) = e^{i(a-b)\phi}\Psi_{a,b}(\mathbf{x}).$$

It follows that

$$\begin{aligned}
 \mathcal{L}_0^{(a,b)}(R_{(0,e^{i\phi})}\Phi)(0) &= \int_{\mathbb{R}^2} R_{(0,e^{i\phi})}\Phi(-\mathbf{y})\Psi_{a,b}(\mathbf{y})d^2\mathbf{y} \\
 &= \int_{\mathbb{R}^2} \Phi(\mathcal{A}_{(0,e^{-i\phi})}(-\mathbf{y}))\Psi_{a,b}(\mathbf{y})d^2\mathbf{y} \\
 &= \int_{\mathbb{R}^2} \Phi(-\mathbf{y})\Psi_{a,b}(\mathcal{A}_{(0,e^{i\phi})}\mathbf{y})d^2\mathbf{y} \\
 &= \int_{\mathbb{R}^2} \Phi(-\mathbf{y})e^{i(b-a)\phi}\Psi_{a,b}(\mathbf{y})d^2\mathbf{y} \\
 &= e^{i(b-a)\phi}\mathcal{L}_0^{(a,b)}(\Phi)(0),
 \end{aligned}$$

i.e.,  $\mathcal{L}_0^{(a,b)}(\Phi)(0)$  transforms under rotations  $e^{i\phi} \in \text{SO}(2)$  as a character (4.22) with  $\xi = b - a$ .

Now consider the map  $\mathcal{L}_{\text{mult}}$ . Since each component in the decomposition (4.50) transforms as a character with  $\xi = b - a$ , the construction of  $\mathcal{L}_{\text{mult}}$  in Sect. 4.2 (based on the procedure of generating invariant polynomials described in Sect. 4.1.4) guarantees that  $\mathcal{L}_{\text{mult}}$  computes a function invariant with respect to  $\text{SO}(2)$ , thus proving Eq. (4.49):

$$\widehat{f}_0(R_{(0,\theta)}\Phi)(0) = \mathcal{L}_{\text{mult}}(\mathcal{L}_{\text{conv}}(R_{(0,\theta)}\Phi)(0)) = \mathcal{L}_{\text{mult}}(\mathcal{L}_{\text{conv}}(\Phi)(0)) = \widehat{f}_0(\Phi)(0).$$

This completes the proof of the necessity part.

**Sufficiency** (*a continuous SE(2)-equivariant map  $f : V \rightarrow U$  can be approximated by charge-conserving convnets*).

Given a continuous SE(2)-equivariant  $f : V \rightarrow U$ , a compact set  $K \subset V$  and positive numbers  $\epsilon, \Lambda_0$ , we need to construct a multi-resolution charge-conserving convnet  $\widehat{f} = (\widehat{f}_\lambda)$  with  $\Lambda > \Lambda_0$  and the property  $\sup_{\Phi \in K} \|\widehat{f}_\lambda(\Phi) - f(\Phi)\| \leq \epsilon$  for all sufficiently small  $\lambda$ . We construct the desired convnet by performing a series of reductions of this approximation problem.

**1. Smoothing.** For any  $\epsilon_1 > 0$ , consider the smoothed map  $\widetilde{f}_{\epsilon_1} : V \rightarrow U$  defined by

$$\widetilde{f}_{\epsilon_1}(\Phi) = f(\Phi) * g_{\epsilon_1}, \tag{4.51}$$

where

$$g_{\epsilon_1}(\mathbf{x}) = \frac{1}{2\pi\epsilon_1} e^{-|\mathbf{x}|^2/(2\epsilon_1)}.$$

The map  $\widetilde{f}_{\epsilon_1}$  is continuous and SE(2)-equivariant, as a composition of two continuous and SE(2)-equivariant maps. We can choose  $\epsilon_1$  small enough so that for all  $\Phi \in K$

$$\|\widetilde{f}_{\epsilon_1}(\Phi) - f(\Phi)\| \leq \frac{\epsilon}{10}. \tag{4.52}$$



The problem of approximating  $f$  then reduces to the problem of approximating maps  $\tilde{f}_{\epsilon_1}$  of the form (4.51).

**2. Spatial cutoff.** We can choose  $\Lambda$  sufficiently large so that for all  $\Phi \in K$

$$\|P_\Lambda \tilde{f}_{\epsilon_1}(\Phi) - \tilde{f}_{\epsilon_1}(\Phi)\| < \frac{\epsilon}{10}. \tag{4.53}$$

We can do this because  $\tilde{f}_{\epsilon_1}(K)$  is compact, as an image of a compact set under a continuous map, and because  $P_\Lambda$  converge strongly to the identity as  $\Lambda \rightarrow +\infty$ . Thus, we only need to approximate output signals  $\tilde{f}_{\epsilon_1}(\Phi)$  on the domain  $[-\Lambda, \Lambda]^2$ .

**3. Output localization.** Define the map  $\tilde{f}_{\epsilon_1, \text{loc}} : V \rightarrow \mathbb{R}$  by

$$\tilde{f}_{\epsilon_1, \text{loc}}(\Phi) = \tilde{f}_{\epsilon_1}(\Phi)(0) = \langle g_{\epsilon_1}, f(\Phi) \rangle_{L^2(\mathbb{R}^2)}. \tag{4.54}$$

Since both  $g_{\epsilon_1}, f(\Phi) \in L^2(\mathbb{R}^2)$ , the map  $\tilde{f}_{\epsilon_1, \text{loc}}$  is well-defined, and it is continuous since  $f$  is continuous.

By translation equivariance of  $f$  and hence  $\tilde{f}_{\epsilon_1}$ , the map  $\tilde{f}_{\epsilon_1}$  can be recovered from  $\tilde{f}_{\epsilon_1, \text{loc}}$  by

$$\tilde{f}_{\epsilon_1}(\Phi)(\mathbf{x}) = \tilde{f}_{\epsilon_1}(R_{(-\mathbf{x}, 1)}\Phi)(0) = \tilde{f}_{\epsilon_1, \text{loc}}(R_{(-\mathbf{x}, 1)}\Phi). \tag{4.55}$$

By the  $\text{SO}(2)$ -equivariance of  $\tilde{f}_{\epsilon_1}$ , the map  $\tilde{f}_{\epsilon_1, \text{loc}}$  is  $\text{SO}(2)$ -invariant.

**4. Nested finite-dimensional  $\text{SO}(2)$ -modules  $V_\zeta$ .** For any nonnegative integer  $a, b$  consider again the signal  $\Psi_{a,b}$  introduced in Eq. (4.20). For any  $\zeta = 1, 2, \dots$ , consider the subspace  $V_\zeta \subset V$  spanned by the vectors  $\text{Re}(\Psi_{a,b})$  and  $\text{Im}(\Psi_{a,b})$  with  $a + b \leq \zeta$ . These vectors form a total system in  $V$  if  $a, b$  take arbitrary nonnegative integer values. Accordingly, if we denote by  $P_{V_\zeta}$  the orthogonal projection to  $V_\zeta$  in  $V$ , then the operators  $P_{V_\zeta}$  converge strongly to the identity as  $\zeta \rightarrow \infty$ .

The subspace  $V_\zeta$  is a real finite-dimensional  $\text{SO}(2)$ -module. As discussed in Sect. 4.1.4, it is convenient to think of such modules as consisting of complex conjugate irreducible representations under constraint (4.29). The complex extension of the real module  $V_\zeta$  is spanned by signals  $\{\Psi_{a,b}\}_{a+b \leq \zeta}$ , so that  $\Psi_{a,b}$  and  $\Psi_{b,a}$  form a complex conjugate pair for  $a \neq b$  (if  $a = b$ , then  $\Psi_{a,b}$  is real). The natural representation (4.1) of  $\text{SO}(2)$  transforms the signal  $\Psi_{a,b}$  as a character (4.22) with  $\xi = a - b$  (in the language of Sect. 4.1.2,  $\Psi_{a,b}$  has local charge  $\eta = b - a$  w.r.t.  $\mathbf{x} = 0$ ):

$$R_{(0, e^{i\phi})}\Psi_{a,b}(\mathbf{x}) = \Psi_{a,b}(\mathcal{A}_{(0, e^{-i\phi})}\mathbf{x}) = e^{i(a-b)\phi}\Psi_{a,b}(\mathbf{x}). \tag{4.56}$$

The action of  $\text{SO}(2)$  on the real signals  $\text{Re}(\Psi_{a,b})$  and  $\text{Im}(\Psi_{a,b})$  can be related to its action on  $\Psi_{a,b}$  and  $\Psi_{b,a}$  as in Eqs.(4.27), (4.28).

**5. Restriction to  $V_\zeta$ .** Let  $\tilde{f}_{\epsilon_1, \text{loc}, \zeta} : V_\zeta \rightarrow \mathbb{R}$  be the restriction of the map  $\tilde{f}_{\epsilon_1, \text{loc}}$  defined in Eq. (4.54) to the subspace  $V_\zeta$ :

$$\tilde{f}_{\epsilon_1, \text{loc}, \zeta} = \tilde{f}_{\epsilon_1, \text{loc}}|_{V_\zeta}. \tag{4.57}$$

Consider the map  $\tilde{f}_{\epsilon_1, \zeta} : V \rightarrow U$  defined by projecting to  $V_\zeta$  and translating the map  $\tilde{f}_{\epsilon_1, \text{loc}, \zeta}$  to points  $\mathbf{x} \in [-\Lambda, \Lambda]^2$  like in the reconstruction formula (4.55):

$$\tilde{f}_{\epsilon_1, \zeta}(\Phi)(\mathbf{x}) = \begin{cases} \tilde{f}_{\epsilon_1, \text{loc}, \zeta}(P_{V_\zeta} R_{(-\mathbf{x}, 1)} \Phi), & \mathbf{x} \in [-\Lambda, \Lambda]^2 \\ 0, & \text{otherwise.} \end{cases} \tag{4.58}$$

We claim that if  $\zeta$  is sufficiently large then for all  $\Phi \in K$

$$\|\tilde{f}_{\epsilon_1, \zeta}(\Phi) - P_\Lambda \tilde{f}_{\epsilon_1}(\Phi)\| < \frac{\epsilon}{10}. \tag{4.59}$$

Indeed,

$$\|\tilde{f}_{\epsilon_1, \zeta}(\Phi) - P_\Lambda \tilde{f}_{\epsilon_1}(\Phi)\| \leq 2\Lambda \sup_{\Phi_1 \in K_1} |\tilde{f}_{\epsilon_1, \text{loc}}(P_{V_\zeta} \Phi_1) - \tilde{f}_{\epsilon_1, \text{loc}}(\Phi_1)|, \tag{4.60}$$

where

$$K_1 = \{R_{(-\mathbf{x}, 1)} \Phi | (\mathbf{x}, \Phi) \in [-\Lambda, \Lambda]^2 \times K\} \subset V. \tag{4.61}$$

The set  $K_1$  is compact, by compactness of  $K$  and strong continuity of  $R$ . Then, by compactness of  $K_1$ , strong convergence  $P_{V_\zeta} \Phi_1 \xrightarrow{\zeta \rightarrow \infty} \Phi_1$  and continuity of  $\tilde{f}_{\epsilon_1, \text{loc}}$ , the r.h.s. of (4.60) becomes arbitrarily small as  $\zeta \rightarrow \infty$ .

It follows from (4.59) that the problem of approximating  $f$  reduces to approximating the map  $\tilde{f}_{\epsilon_1, \zeta}$  for a fixed finite  $\zeta$ .

**6. Polynomial approximation.** The map  $\tilde{f}_{\epsilon_1, \text{loc}, \zeta} : V_\zeta \rightarrow \mathbb{R}$  defined in (4.57) is a continuous  $\text{SO}(2)$ -invariant map on the  $\text{SO}(2)$ -module  $V_\zeta$ . By Lemma 4.2, such a map can be approximated by invariant polynomials. Let  $K_1 \subset V$  be the compact set defined in Eq. (4.61). Note that  $P_{V_\zeta} K_1$  is then a compact subset of  $V_\zeta$ . Let  $\widehat{f}_{\text{loc}} : V_\zeta \rightarrow \mathbb{R}$  be an  $\text{SO}(2)$ -invariant polynomial such that for all  $\Phi_2 \in P_{V_\zeta} K_1$

$$|\widehat{f}_{\text{loc}}(\Phi_2) - \tilde{f}_{\epsilon_1, \text{loc}, \zeta}(\Phi_2)| \leq \frac{\epsilon}{10 \cdot 2\Lambda}. \tag{4.62}$$

Consider now the map  $\widehat{f}_0 : V \rightarrow U$  defined by

$$\widehat{f}_0(\Phi)(\mathbf{x}) = \begin{cases} \widehat{f}_{\text{loc}}(P_{V_\zeta} R_{(-\mathbf{x}, 1)} \Phi), & \mathbf{x} \in [-\Lambda, \Lambda]^2, \\ 0, & \text{otherwise.} \end{cases} \tag{4.63}$$

Using Eqs.(4.58) and (4.62), we have for all  $\mathbf{x} \in [-\Lambda, \Lambda]^2$  and  $\Phi_2 \in P_{V_\zeta} K_1$

$$|\widehat{f}_0(\Phi_2)(\mathbf{x}) - \tilde{f}_{\epsilon_1, \zeta}(\Phi_2)(\mathbf{x})| \leq \frac{\epsilon}{10 \cdot 2\Lambda}$$

and hence for all  $\Phi \in K$

$$\|\widehat{f}_0(\Phi) - \widetilde{f}_{\epsilon_1, \zeta}(\Phi)\| < \frac{\epsilon}{10}. \tag{4.64}$$

**7. Identification of convnet with  $\lambda = 0$ .** We show now that the map  $\widehat{f}_0$  given in (4.63) can be written as the scaling limit ( $\lambda \rightarrow 0$ ) of a multi-resolution charge-conserving convnet.

First note that the projector  $P_{V_\zeta}$  can be written as

$$P_{V_\zeta} \Phi = \sum_{a,b:a+b \leq \zeta} \langle \Psi'_{a,b}, \Phi \rangle \Psi_{a,b},$$

where  $\Psi'_{a,b}$  is the basis in  $V_\zeta$  dual to the basis  $\Psi_{a,b}$ . Let  $V_{\zeta, \xi}$  denote the isotypic component in  $V_\zeta$  spanned by vectors  $\Psi_{a,b}$  with  $a - b = \xi$ . By Eq. (4.56), this notation is consistent with the notation of Sect. 4.1.4 where the number  $\xi$  is used to specify the characters (4.22). By unitarity of the representation  $R$ , different isotypic components are mutually orthogonal, so  $\Psi'_{a,b} \in V_{\zeta, a-b}$  and we can expand

$$\Psi'_{a,b} = \sum_{\substack{0 \leq a', b' \leq \zeta \\ a' - b' = a - b}} c_{a,b,a',b'} \Psi_{a',b'}$$

with some coefficients  $c_{a,b,a',b'}$ . Then we can write

$$\begin{aligned} P_{V_\zeta} R_{(-x,1)} \Phi &= \sum_{a,b:a+b \leq \zeta} \langle \Psi'_{a,b}, R_{(-x,1)} \Phi \rangle \Psi_{a,b} \\ &= \sum_{\xi=-\zeta}^{\zeta} \sum_{\substack{a,b:a+b \leq \zeta \\ a-b=\xi}} \sum_{\substack{a',b':a'+b' \leq \zeta \\ a'-b'=\xi}} \overline{c_{a,b,a',b'}} \langle \Psi_{a',b'}, R_{(-x,1)} \Phi \rangle \Psi_{a,b} \\ &= \sum_{\xi=-\zeta}^{\zeta} \sum_{\substack{a,b:a+b \leq \zeta \\ a-b=\xi}} \sum_{\substack{a',b':a'+b' \leq \zeta \\ a'-b'=\xi}} \overline{c_{a,b,a',b'}} \left( \int_{\mathbb{R}^2} \Phi(\mathbf{x} + \mathbf{y}) \overline{\Psi_{a',b'}(\mathbf{y})} d^2 \mathbf{y} \right) \Psi_{a,b} \\ &= \sum_{\xi=-\zeta}^{\zeta} \sum_{\substack{a,b:a+b \leq \zeta \\ a-b=\xi}} \sum_{\substack{a',b':a'+b' \leq \zeta \\ a'-b'=\xi}} \overline{c_{a,b,a',b'}} (-1)^{a'+b'} \left( \int_{\mathbb{R}^2} \Phi(\mathbf{x} - \mathbf{y}) \Psi_{b',a'}(\mathbf{y}) d^2 \mathbf{y} \right) \Psi_{a,b} \\ &= \sum_{\xi=-\zeta}^{\zeta} \sum_{\substack{a,b:a+b \leq \zeta \\ a-b=\xi}} \sum_{\substack{a',b':a'+b' \leq \zeta \\ a'-b'=\xi}} \overline{c_{a,b,a',b'}} (-1)^{a'+b'} (\mathcal{L}_0^{(b',a')} \Phi(\mathbf{x})) \Psi_{a,b}, \end{aligned} \tag{4.65}$$

where in the penultimate step we used the identity  $\overline{\Psi_{a',b'}(\mathbf{y})} = \Psi_{(b',a')}(\mathbf{y}) = (-1)^{a'+b'} \Psi_{(b',a')}(-\mathbf{y})$ , and in the last step we used definition (4.21) of  $\mathcal{L}_0^{(a,b)}$ .

We can now interpret the map  $\widehat{f}_0$  given by (4.63) as the  $\lambda \rightarrow 0$  limit of a convnet of Sect. 4.2 in the following way.

First, by the above expansion, the part  $P_{V_\zeta} R_{(-\mathbf{x}, 1)}$  of the map  $\widehat{f}_0$  computes various convolutions  $\mathcal{L}_0^{(b', a')} \Phi$  with  $a' + b' \leq \zeta$ ,  $a' - b' = \xi$ —this corresponds to the  $\lambda \rightarrow 0$  limit of smoothing and differentiation layers of Sect. 4.2 with  $T_{\text{diff}} = \zeta$ . The global charge parameter  $\mu$  appearing in the decomposition (4.32) of the target spaces  $W_{\text{diff}, t}$  of differentiation layers corresponds to  $-\xi (= b' - a')$  in the above formula, while the degree  $s$  corresponds to  $a' + b'$ . The vectors  $\Psi_{a, b}$  with  $a - b = \xi$  over which we expand in (4.65) serve as a particular basis in the  $\mu = -\xi$  component of  $W_{\text{diff}, T_{\text{diff}}}$ .

Now, the invariant polynomial  $\widehat{f}_{\text{loc}}$  appearing in (4.63) can be expressed as a polynomial in the variables associated with the isotypic components  $V_{\zeta, \xi}$ . These components are spanned by the vectors  $\Psi_{a, b}$  with  $a - b = \xi$ . By Eq. (4.65),  $\widehat{f}_{\text{loc}}(P_{V_\zeta} R_{(-\mathbf{x}, 1)} \Phi)$  can then be viewed as an invariant polynomial in the variables  $\mathcal{L}_0^{(b', a')} \Phi(\mathbf{x})$  that correspond to the isotypic components  $V_{\zeta, \xi}$  with  $\xi = a' - b'$ . As shown in Sect. 4.1.4, this invariant polynomial can then be generated by the layerwise multiplication procedure (4.26) starting from the initial variables  $\mathcal{L}_0^{(b', a')} \Phi(\mathbf{x})$ . This procedure is reproduced in the definition (4.36), (4.37) of convnet multiplication layers. (The charge-conservation constraints are expressed in Eqs.(4.36), (4.37) in terms of  $\mu$  rather than  $\xi$ , but  $\mu = -\xi$ , and the constraints are invariant with respect to changing the sign of all  $\mu$ 's.) Thus, if the number  $T_{\text{mult}}$  of multiplication layers and the dimensions  $d_{\text{mult}}$  of these layers are sufficiently large, then one can arrange the weights in these layers so as to exactly give the map  $\Phi \mapsto \widehat{f}_{\text{loc}}(P_{V_\zeta} R_{(-\mathbf{x}, 1)} \Phi)$ .

**8. Approximation by convnets with  $\lambda > 0$ .** It remains to show that the scaling limit  $\widehat{f}_0$  is approximated by the  $\lambda > 0$  convnets  $\widehat{f}_\lambda$  in the sense that if  $\lambda$  is sufficiently small then for all  $\Phi \in K$

$$\|\widehat{f}_0(\Phi) - \widehat{f}_\lambda(\Phi)\| < \frac{\epsilon}{10}. \tag{4.66}$$

We have already shown earlier in Eq. (4.42) that for any  $\Phi \in V$  the signals  $\widehat{f}_\lambda(\Phi)$  converge to  $\widehat{f}_0(\Phi)$  in the  $L^2([-\Lambda, \Lambda]^2)$  sense. In fact, Lemma 4.1 implies that this convergence is uniform on any compact set  $K \subset V$ , which proves Eq. (4.66).

Summarizing all the above steps, we have constructed a multi-resolution charge-conserving convnet  $\widehat{f}_\lambda$  such that, by the inequalities (4.52), (4.53), (4.59), (4.64) and (4.66), we have  $\sup_{\Phi \in K} \|\widehat{f}_\lambda(\Phi) - f(\Phi)\| \leq \epsilon$  for all sufficiently small  $\lambda$ . This completes the proof of the sufficiency part.  $\square$

### 5 Discussion

We summarize and discuss the obtained results, and indicate potential directions of further research.

In Sect. 2 we considered approximation of maps defined on finite-dimensional spaces and described universal and exactly invariant/equivariant extensions of the usual shallow neural network (Propositions 2.3, 2.4). These extensions are obtained by adding to the network a special polynomial layer. This construction can be seen as an alternative to the symmetrization of the network (similarly to how constructing

symmetric polynomials as functions of elementary symmetric polynomials is an alternative to symmetrizing non-symmetric polynomials). A drawback (inherited from the theory of invariant polynomials) of this construction is that it requires us to know appropriate sets of generating polynomial invariants/equivariants, which is difficult in practice. This difficulty can be ameliorated using polarization if the modules in question are decomposed into multiple copies of a few basic modules (Proposition 2.5, 2.7), but this approach still may be too complicated in general for practical applications.

Nevertheless, in the case of the symmetric group  $S_N$  we have derived an explicit complete  $S_N$ -invariant modification of the usual shallow neural network (Theorem 2.4). While complete and exactly  $S_N$ -invariant, this modification does not involve symmetrization over  $S_N$ . With its relatively small computational complexity, this modification thus presents a viable alternative to the symmetrization-based approach.

One can expect that further progress in the design of invariant/equivariant models may be achieved by using more advanced general constructions from the representation and invariant theories. In particular, in Sect. 2 we have not considered *products* of representations, but later in Sect. 4 we essentially use them in the abelian  $SO(2)$  setting when defining multiplication layers in “charge-conserving convnet”.

In Sect. 3 we considered approximations of maps defined on the space  $V = L^2(\mathbb{R}^v, \mathbb{R}^{dv})$  of  $d_V$ -component signals on  $\mathbb{R}^v$ . The crucial feature of this setting is the infinite-dimensionality of the space  $V$ , which requires us to reconsider the notion of approximation. Inspired by classical finite-dimensional results [32], our approach in Sect. 3 was to assume that a map  $f$  is defined on the whole  $L^2(\mathbb{R}^v, \mathbb{R}^{dv})$  as a map  $f : L^2(\mathbb{R}^v, \mathbb{R}^{dv}) \rightarrow L^2(\mathbb{R}^v, \mathbb{R}^{dv})$  or  $f : L^2(\mathbb{R}^v, \mathbb{R}^{dv}) \rightarrow \mathbb{R}$ , and consider its approximation by finite models  $\tilde{f}$  in a weak sense of comparison on compact subsets of  $V$  (see Definitions 3.2 and 3.4). This approach has allowed us to prove reasonable universal approximation properties of standard convnets. Specifically, in Theorem 3.1 we prove that a map  $f : L^2(\mathbb{R}^v, \mathbb{R}^{dv}) \rightarrow L^2(\mathbb{R}^v, \mathbb{R}^{dv})$  can be approximated by convnets without pooling if and only if  $f$  is norm-continuous and  $\mathbb{R}^v$ -equivariant. In Theorem 3.2 we prove that a map  $f : L^2(\mathbb{R}^v, \mathbb{R}^{dv}) \rightarrow \mathbb{R}$  can be approximated by convnets with downsampling if and only if  $f$  is norm-continuous.

In applications involving convnets (e.g., image recognition or segmentation), the approximated maps  $f$  are considered only on small subsets of the full space  $V$ . Compact (or, more generally, precompact) subsets have properties that seem to make them a reasonable general abstraction for such subsets. In particular, a subset  $K \subset V$  is precompact if, for example, it results from a continuous generative process involving finitely many bounded parameters; or if  $K$  is a finite union of precompact subsets; or if for any  $\epsilon > 0$  the set  $K$  can be covered by finitely many  $\epsilon$ -balls. From this perspective, it seems reasonable to consider restrictions of maps  $f$  to compact sets, as we did in our weak notion of approximation in Sect. 3. At the same time, it would be interesting to refine the notion of model convergence by considering the structure of the sets  $K$  in more detail and relate it quantitatively to the approximation accuracy (in particular, paving the way to computing approximation rates).

In Sect. 4 we consider the task of constructing finite universal approximators for maps  $f : L^2(\mathbb{R}^2, \mathbb{R}^{dv}) \rightarrow L^2(\mathbb{R}^2, \mathbb{R}^{dv})$  equivariant with respect to the group  $SE(2)$  of two-dimensional rigid planar motions. We introduce a particular convnet-like model—“charge-conserving convnet”—solving this task. We extend the topological framework

of Sect. 3 to rigorously formulate the properties of equivariance and completeness to be proved. Our main result, Theorem 4.1, shows that a map  $f : L^2(\mathbb{R}^2, \mathbb{R}^{d\nu}) \rightarrow L^2(\mathbb{R}^2, \mathbb{R}^{d\nu})$  can be approximated in the small-scale limit by finite charge-conserving convnets if and only if  $f$  is norm-continuous and SE(2)-equivariant.

The construction of this convnet is based on splitting the feature space into isotopic components characterized by a particular representation of the group SO(2) of proper 2D rotations. The information flow in the model is constrained by what can be interpreted as “charge conservation” (hence the name of the model). The model is essentially polynomial, only including elementary arithmetic operations (+, −, \*) arranged so as to satisfy these constraints but otherwise achieve full expressivity.

While in Sects. 3, 4 we have constructed intrinsically  $\mathbb{R}^\nu$ - and (for  $\nu = 2$ ) SO(2)-equivariant and complete approximators for maps  $f : L^2(\mathbb{R}^\nu, \mathbb{R}^{d\nu}) \rightarrow L^2(\mathbb{R}^\nu, \mathbb{R}^{d\nu})$ , we have not been able to similarly construct intrinsically  $\mathbb{R}^\nu$ -invariant approximators for maps  $f : L^2(\mathbb{R}^\nu, \mathbb{R}^{d\nu}) \rightarrow \mathbb{R}$ . As noted in Sect. 3 and confirmed by Theorem 3.2, if we simply include pooling in the convnet, it completely destroys the  $\mathbb{R}^\nu$ -invariance in our continuum limit. It would be interesting to further explore this issue.

The convnets considered in Sect. 3 have a rather conventional structure as sequences of linear convolutional layers equipped with a nonlinear activation function [12]. In contrast, the charge-conserving convnets of Sect. 4 have a special and somewhat artificial structure (three groups of layers of which the first two are linear and commuting; no arbitrary nonlinearities). This structure was essential for our proof of the main Theorem 4.1, since these assumptions on the model allowed us to prove that the model is both SE(2)-equivariant and complete. It would be interesting to extend this theorem to more general approximation models.

**Acknowledgements** The author thanks the anonymous reviewer for several helpful suggestions.

## A Proof of Lemma 4.1

The proof is a slight modification of the standard proof of Central Limit Theorem via Fourier transform (the CLT can be directly used to prove the lemma in the case  $a = b = 0$  when  $\mathcal{L}_\lambda^{(a,b)}$  only includes diffusion factors).

To simplify notation, assume without loss of generality that  $d_V = 1$  (in the general case the proof is essentially identical). We will use the appropriately discretized version of the Fourier transform (i.e., the Fourier series expansion). Given a discretized signal  $\Phi : (\lambda\mathbb{Z})^2 \rightarrow \mathbb{C}$ , we define  $\mathcal{F}_\lambda \Phi$  as a function on  $[-\frac{\pi}{\lambda}, \frac{\pi}{\lambda}]^2$  by

$$\mathcal{F}_\lambda \Phi(\mathbf{p}) = \frac{\lambda^2}{2\pi} \sum_{\gamma \in (\lambda\mathbb{Z})^2} \Phi(\gamma) e^{-i\mathbf{p}\cdot\gamma}.$$

Then,  $\mathcal{F}_\lambda : L^2((\lambda\mathbb{Z})^2, \mathbb{C}) \rightarrow L^2([-\frac{\pi}{\lambda}, \frac{\pi}{\lambda}]^2, \mathbb{C})$  is a unitary isomorphism, assuming that the scalar product in the input space is defined by  $\langle \Phi, \Psi \rangle = \lambda^2 \sum_{\gamma \in (\lambda\mathbb{Z})^2} \overline{\Phi(\gamma)} \Psi(\gamma)$  and in the output space by  $\langle \Phi, \Psi \rangle = \int_{[-\frac{\pi}{\lambda}, \frac{\pi}{\lambda}]^2} \overline{\Phi(\mathbf{p})} \Psi(\mathbf{p}) d^2\mathbf{p}$ . Let  $P_\lambda$  be the discretization projector (3.6). It is easy to check that  $\mathcal{F}_\lambda P_\lambda$  strongly converges to the standard

Fourier transform as  $\lambda \rightarrow 0$  :

$$\lim_{\lambda \rightarrow 0} \mathcal{F}_\lambda P_\lambda \Phi = \mathcal{F}_0 \Phi, \quad \Phi \in L^2(\mathbb{R}^2, \mathbb{C}),$$

where

$$\mathcal{F}_0 \Phi(\mathbf{p}) = \frac{1}{2\pi} \int_{\mathbb{R}^2} \Phi(\gamma) e^{-i\mathbf{p} \cdot \gamma} d^2 \gamma$$

and where we naturally embed  $L^2([-\frac{\pi}{\lambda}, \frac{\pi}{\lambda}]^2, \mathbb{C}) \subset L^2(\mathbb{R}^2, \mathbb{C})$ . Conversely, let  $P'_\lambda$  denote the orthogonal projection onto the subspace  $L^2([-\frac{\pi}{\lambda}, \frac{\pi}{\lambda}]^2, \mathbb{C})$  in  $L^2(\mathbb{R}^2, \mathbb{C})$  :

$$P'_\lambda : \Phi \mapsto \Phi|_{[-\frac{\pi}{\lambda}, \frac{\pi}{\lambda}]^2}. \tag{A.1}$$

Then

$$\lim_{\lambda \rightarrow 0} \mathcal{F}_\lambda^{-1} P'_\lambda \Phi = \mathcal{F}_0^{-1} \Phi, \quad \Phi \in L^2(\mathbb{R}^2, \mathbb{C}). \tag{A.2}$$

Fourier transform gives us the spectral representation of the discrete differential operators (4.15), (4.16), (4.17) as operators of multiplication by function:

$$\begin{aligned} \mathcal{F}_\lambda \partial_z^{(\lambda)} \Phi &= \Psi_{\partial_z^{(\lambda)}} \cdot \mathcal{F}_\lambda \Phi, \\ \mathcal{F}_\lambda \partial_{\bar{z}}^{(\lambda)} \Phi &= \Psi_{\partial_{\bar{z}}^{(\lambda)}} \cdot \mathcal{F}_\lambda \Phi, \\ \mathcal{F}_\lambda \Delta^{(\lambda)} \Phi &= \Psi_{\Delta^{(\lambda)}} \cdot \mathcal{F}_\lambda \Phi, \end{aligned}$$

where, denoting  $\mathbf{p} = (p_x, p_y)$ ,

$$\begin{aligned} \Psi_{\partial_z^{(\lambda)}}(p_x, p_y) &= \frac{i}{2\lambda} (\sin \lambda p_x - i \sin \lambda p_y), \\ \Psi_{\partial_{\bar{z}}^{(\lambda)}}(p_x, p_y) &= \frac{i}{2\lambda} (\sin \lambda p_x + i \sin \lambda p_y), \\ \Psi_{\Delta^{(\lambda)}}(p_x, p_y) &= -\frac{4}{\lambda^2} \left( \sin^2 \frac{\lambda p_x}{2} + \sin^2 \frac{\lambda p_y}{2} \right). \end{aligned}$$

The operator  $\mathcal{L}_\lambda^{(a,b)}$  defined in (4.19) can then be written as

$$\mathcal{F}_\lambda \mathcal{L}_\lambda^{(a,b)} \Phi = \Psi_{\mathcal{L}_\lambda^{(a,b)}} \cdot \mathcal{F}_\lambda P_\lambda \Phi,$$

where the function  $\Psi_{\mathcal{L}_\lambda^{(a,b)}}$  is given by

$$\Psi_{\mathcal{L}_\lambda^{(a,b)}} = (\Psi_{\partial_z^{(\lambda)}})^a (\Psi_{\partial_{\bar{z}}^{(\lambda)}})^b \left( 1 + \frac{\lambda^2}{8} \Psi_{\Delta^{(\lambda)}} \right)^{\lceil 4/\lambda^2 \rceil}.$$

We can then write  $\mathcal{L}_\lambda^{(a,b)}\Phi$  as a convolution of  $P_\lambda\Phi$  with the kernel

$$\Psi_{a,b}^{(\lambda)} = \frac{1}{2\pi} \mathcal{F}_\lambda^{-1} \Psi_{\mathcal{L}_\lambda^{(a,b)}}$$

on the grid  $(\lambda\mathbb{Z})^2$  :

$$\mathcal{L}_\lambda^{(a,b)}\Phi(\gamma) = \lambda^2 \sum_{\theta \in (\lambda\mathbb{Z})^2} P_\lambda\Phi(\gamma - \theta) \Psi_{a,b}^{(\lambda)}(\theta), \quad \gamma \in (\lambda\mathbb{Z})^2. \tag{A.3}$$

Now consider the operator  $\mathcal{L}_0^{(a,b)}$  defined in (4.21). At each  $\mathbf{x} \in \mathbb{R}^2$ , the value  $\mathcal{L}_0^{(a,b)}\Phi(\mathbf{x})$  can be written as a scalar product:

$$\mathcal{L}_0^{(a,b)}\Phi(\mathbf{x}) = \int_{\mathbb{R}^2} \Phi(\mathbf{x} - \mathbf{y}) \Psi_{a,b}(\mathbf{y}) d^2\mathbf{y} = \langle R_{-\mathbf{x}}\tilde{\Phi}, \Psi_{a,b} \rangle_{L^2(\mathbb{R}^2)}, \tag{A.4}$$

where  $\tilde{\Phi}(\mathbf{x}) = \overline{\Phi(-\mathbf{x})}$ ,  $\Psi_{a,b}$  is defined by (4.20), and  $R_{\mathbf{x}}$  is our standard representation of the group  $\mathbb{R}^2$ ,  $R_{\mathbf{x}}\Phi(\mathbf{y}) = \Phi(\mathbf{y} - \mathbf{x})$ . For  $\lambda > 0$ , we can write  $\mathcal{L}_\lambda^{(a,b)}\Phi(\mathbf{x})$  in a similar form. Indeed, using (A.3) and naturally extending the discretized signal  $\Psi_{a,b}^{(\lambda)}$  to the whole  $\mathbb{R}^2$ , we have

$$\mathcal{L}_\lambda^{(a,b)}\Phi(\gamma) = \int_{\mathbb{R}^2} \Phi(\gamma - \mathbf{y}) \Psi_{a,b}^{(\lambda)}(\mathbf{y}) d^2\mathbf{y} = \langle R_{-\gamma}\tilde{\Phi}, \Psi_{a,b}^{(\lambda)} \rangle_{L^2(\mathbb{R}^2)}.$$

Then, for any  $\mathbf{x} \in \mathbb{R}^2$  we can write

$$\mathcal{L}_\lambda^{(a,b)}\Phi(\mathbf{x}) = \langle R_{-\mathbf{x} + \delta\mathbf{x}}\tilde{\Phi}, \Psi_{a,b}^{(\lambda)} \rangle_{L^2(\mathbb{R}^2)}, \tag{A.5}$$

where  $-\mathbf{x} + \delta\mathbf{x}$  is the point of the grid  $(\lambda\mathbb{Z})^2$  nearest to  $-\mathbf{x}$ .

Now consider the formulas (A.4), (A.5) and observe that, by Cauchy-Schwarz inequality and since  $R$  is norm-preserving, to prove statement 1) of the lemma we only need to show that the functions  $\Psi_{a,b}$ ,  $\Psi_{a,b}^{(\lambda)}$  have uniformly bounded  $L^2$ -norms. For  $\lambda > 0$  we have

$$\begin{aligned} \|\Psi_{a,b}^{(\lambda)}\|_{L^2(\mathbb{R}^2)}^2 &= \left\| \frac{1}{2\pi} \mathcal{F}_\lambda^{-1} \Psi_{\mathcal{L}_\lambda^{(a,b)}} \right\|_{L^2(\mathbb{R}^2)}^2 \\ &= \frac{1}{4\pi^2} \|\Psi_{\mathcal{L}_\lambda^{(a,b)}}\|_{L^2(\mathbb{R}^2)}^2 \\ &= \frac{1}{4\pi^2} \|(\Psi_{\partial_z^{(\lambda)}})^a (\Psi_{\partial_{\bar{z}}^{(\lambda)}})^b \left(1 + \frac{\lambda^2}{8} \Psi_{\Delta^{(\lambda)}}\right)^{\lceil 4/\lambda^2 \rceil}\|_{L^2(\mathbb{R}^2)}^2 \\ &\leq \frac{1}{4\pi^2} \int_{-\pi/\lambda}^{\pi/\lambda} \int_{-\pi/\lambda}^{\pi/\lambda} \left(\frac{|p_x| + |p_y|}{2}\right)^{2(a+b)} \\ &\quad \exp\left(-\lceil 4/\lambda^2 \rceil \left(\sin^2 \frac{\lambda p_x}{2} + \sin^2 \frac{\lambda p_y}{2}\right)\right) dp_x dp_y \end{aligned}$$



$$\begin{aligned} &\leq \frac{1}{4\pi^2} \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \left(\frac{|p_x|+|p_y|}{2}\right)^{2(a+b)} \exp\left(-\frac{4}{\pi^2}(p_x^2 + p_y^2)\right) dp_x dp_y \\ &< \infty, \end{aligned} \tag{A.6}$$

where we used the inequalities

$$\begin{aligned} |\sin t| &\leq |t|, \\ |1 + t| &\leq e^t, \quad t > -1, \\ |\sin t| &\geq \frac{2|t|}{\pi}, \quad t \in \left[-\frac{\pi}{2}, \frac{\pi}{2}\right]. \end{aligned}$$

Expression (A.6) provides a finite bound, uniform in  $\lambda$ , for the squared norms  $\|\Psi_{a,b}^{(\lambda)}\|^2$ . This bound also holds for  $\|\Psi_{a,b}\|^2$ .

Next, observe that to establish the strong convergence in statement 2) of the lemma, it suffices to show that

$$\lim_{\lambda \rightarrow 0} \|\Psi_{a,b}^{(\lambda)} - \Psi_{a,b}\|_{L^2(\mathbb{R}^2)} = 0. \tag{A.7}$$

Indeed, by (A.4), (A.5), we would then have

$$\begin{aligned} \|\mathcal{L}_\lambda^{(a,b)} \Phi - \mathcal{L}_0^{(a,b)} \Phi\|_\infty &= \sup_{\mathbf{x} \in \mathbb{R}^2} |\langle R_{-\mathbf{x}+\delta\mathbf{x}} \tilde{\Phi}, \Psi_{a,b}^{(\lambda)} \rangle - \langle R_{-\mathbf{x}} \tilde{\Phi}, \Psi_{a,b} \rangle| \\ &= \sup_{\mathbf{x} \in \mathbb{R}^2} |\langle R_{-\mathbf{x}}(R_{\delta\mathbf{x}} - 1) \tilde{\Phi}, \Psi_{a,b}^{(\lambda)} \rangle + \langle R_{-\mathbf{x}} \tilde{\Phi}, \Psi_{a,b}^{(\lambda)} - \Psi_{a,b} \rangle| \\ &\leq \sup_{\|\delta\mathbf{x}\| \leq \lambda} \|R_{\delta\mathbf{x}} \tilde{\Phi} - \tilde{\Phi}\|_2 \sup_{\lambda} \|\Psi_{a,b}^{(\lambda)}\|_2 + \|\tilde{\Phi}\|_2 \|\Psi_{a,b}^{(\lambda)} - \Psi_{a,b}\|_2 \\ &\xrightarrow{\lambda \rightarrow 0} 0 \end{aligned}$$

thanks to the unitarity of  $R$ , convergence  $\lim_{\delta\mathbf{x} \rightarrow 0} \|R_{\delta\mathbf{x}} \tilde{\Phi} - \tilde{\Phi}\|_2 = 0$ , uniform boundedness of  $\|\Psi_{a,b}^{(\lambda)}\|_2$  and convergence (A.7).

To establish (A.7), we write

$$\Psi_{a,b}^{(\lambda)} - \Psi_{a,b} = \frac{1}{2\pi} (\mathcal{F}_\lambda^{-1} \Psi_{\mathcal{L}_\lambda^{(a,b)}} - \mathcal{F}_0^{-1} \Psi_{\mathcal{L}_0^{(a,b)}}),$$

where  $\Psi_{\mathcal{L}_0^{(a,b)}} = 2\pi \mathcal{F}_\lambda \Psi_{a,b}$ . By definition (4.20) of  $\Psi_{a,b}$  and standard properties of Fourier transform, the explicit form of the function  $\Psi_{\mathcal{L}_0^{(a,b)}}$  is

$$\Psi_{\mathcal{L}_0^{(a,b)}}(p_x, p_y) = \left(\frac{i(p_x - ip_y)}{2}\right)^a \left(\frac{i(p_x + ip_y)}{2}\right)^b \exp\left(-\frac{p_x^2 + p_y^2}{2}\right).$$

Observe that the function  $\Psi_{\mathcal{L}_0^{(a,b)}}$  is the pointwise limit of the functions  $\Psi_{\mathcal{L}_\lambda^{(a,b)}}$  as  $\lambda \rightarrow 0$ . The functions  $|\Psi_{\mathcal{L}_\lambda^{(a,b)}}|^2$  are bounded uniformly in  $\lambda$  by the integrable function

appearing in the integral (A.6). Therefore we can use the dominated convergence theorem and conclude that

$$\lim_{\lambda \rightarrow 0} \|\Psi_{\mathcal{L}_\lambda^{(a,b)}} - P'_\lambda \Psi_{\mathcal{L}_0^{(a,b)}}\|_2 = 0, \quad (\text{A.8})$$

where  $P'_\lambda$  is the cut-off projector (A.1). We then have

$$\begin{aligned} \|\Psi_{a,b}^{(\lambda)} - \Psi_{a,b}\|_2 &= \frac{1}{2\pi} \|\mathcal{F}_\lambda^{-1} \Psi_{\mathcal{L}_\lambda^{(a,b)}} - \mathcal{F}_0^{-1} \Psi_{\mathcal{L}_0^{(a,b)}}\|_2 \\ &\leq \frac{1}{2\pi} \|\mathcal{F}_\lambda^{-1} (\Psi_{\mathcal{L}_\lambda^{(a,b)}} - P'_\lambda \Psi_{\mathcal{L}_0^{(a,b)}})\|_2 + \frac{1}{2\pi} \|(\mathcal{F}_\lambda^{-1} P'_\lambda - \mathcal{F}_0^{-1}) \Psi_{\mathcal{L}_0^{(a,b)}}\|_2 \\ &\xrightarrow{\lambda \rightarrow 0} 0 \end{aligned}$$

by (A.8) and (A.2). We have thus proved (A.7).

It remains to show that the convergence  $\mathcal{L}_\lambda^{(a,b)} \Phi \rightarrow \mathcal{L}_0^{(a,b)} \Phi$  is uniform on compact sets  $K \subset V$ . This follows by a version of continuity argument. For any  $\epsilon > 0$ , we can choose finitely many  $\Phi_n, n = 1, \dots, N$ , such that for any  $\Phi \in K$  there is some  $\Phi_n$  for which  $\|\Phi - \Phi_n\| < \epsilon$ . Then  $\|\mathcal{L}_\lambda^{(a,b)} \Phi - \mathcal{L}_0^{(a,b)} \Phi\| \leq \|\mathcal{L}_\lambda^{(a,b)} \Phi_n - \mathcal{L}_0^{(a,b)} \Phi_n\| + 2 \sup_{\lambda \geq 0} \|\mathcal{L}_\lambda^{(a,b)}\| \epsilon$ . Since  $\sup_{\lambda \geq 0} \|\mathcal{L}_\lambda^{(a,b)}\| < \infty$  by statement 1) of the lemma, the desired uniform convergence for  $\Phi \in K$  follows from the convergence for  $\Phi_n, n = 1, \dots, N$ .

## References

1. Anselmi, F., Rosasco, L., Poggio, T.: On invariance and selectivity in representation learning. *Inference* **5**(2), 134–158 (2016)
2. Bruna, J., Mallat, S.: Invariant scattering convolution networks. *IEEE Trans. Pattern Anal. Mach. Intell.* **35**(8), 1872–1886 (2013)
3. Burkhart, H., Siggelkow, S.: Invariant features in pattern recognition-fundamentals and applications. *In: Nonlinear Model-Based Image/Video Processing and Analysis*, pp. 269–307 (2001)
4. Cohen, N., Shashua, A.: Convolutional rectifier networks as generalized tensor decompositions. *In: International Conference on Machine Learning*, pp. 955–963 (2016)
5. Cohen, N., Sharir, O., Levine, Y., Tamari, R., Yakira, D., Shashua, A.: Analysis and design of convolutional networks via hierarchical tensor decompositions (2017). [arXiv preprint arXiv:1705.02302](https://arxiv.org/abs/1705.02302)
6. Cohen, T., Welling, M.: Group equivariant convolutional networks. *In: Proceedings of the 33rd International Conference on Machine Learning*, pp. 2990–2999 (2016)
7. Cybenko, G.: Approximation by superpositions of a sigmoidal function. *Math. Control Signals Syst.* **2**(4), 303–314 (1989)
8. Dieleman, S., De Fauw, J., Kavukcuoglu, K.: Exploiting cyclic symmetry in convolutional neural networks. *In: Proceedings of the 33rd International Conference on International Conference on Machine Learning*, vol. 48, pp. 1889–1898 (2016)
9. Esteves, C., Allen-Blanchette, C., Zhou, X., Daniilidis, K.: Polar transformer networks. *In: International Conference on Learning Representations* (2018)
10. Funahashi, K.-I.: On the approximate realization of continuous mappings by neural networks. *Neural Netw.* **2**(3), 183–192 (1989)
11. Gens, R., Domingos, P.M.: Deep symmetry networks. *In: Advances in Neural Information Processing Systems*, pp. 2537–2545 (2014)
12. Goodfellow, I., Bengio, Y.: *Deep Learning*. MIT Press, Cambridge (2016)

13. He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recognition. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pp. 770–778 (2016)
14. Hedlund, G.A.: Endomorphisms and automorphisms of the shift dynamical system. *Theory Comput. Syst.* **3**(4), 320–375 (1969)
15. Henriques, J.F., Vedaldi, A.: Warped convolutions: efficient invariance to spatial transformations. In: International Conference on Machine Learning, pp. 1461–1469 (2017)
16. Hilbert, D.: Über die Theorie der algebraischen Formen. *Mathematische Annalen* **36**(4), 473–534 (1890)
17. Hilbert, D.: Über die vollen Invariantensysteme. *Mathematische Annalen* **42**(3), 313–373 (1893)
18. Hornik, K.: Some new results on neural network approximation. *Neural Netw.* **6**(8), 1069–1072 (1993)
19. Hornik, K., Stinchcombe, M., White, H.: Multilayer feedforward networks are universal approximators. *Neural Netw.* **2**(5), 359–366 (1989)
20. Kondor, R., Trivedi, S.: On the generalization of equivariance and convolution in neural networks to the action of compact groups. In: International Conference on Machine Learning, pp. 2747–2755 (2018)
21. Kraft, H., Procesi, C.: Classical invariant theory, a primer. *Lecture Notes* (2000)
22. le Cun, Y.: Generalization and network design strategies. In: *Connectionism in Perspective*, pp. 143–155 (1989)
23. LeCun, Y., Bengio, Y., Hinton, G.: Deep learning. *Nature* **521**(7553), 436–444 (2015)
24. Leshno, M., Lin, V.Y., Pinkus, A., Schocken, S.: Multilayer feedforward networks with a nonpolynomial activation function can approximate any function. *Neural Netw.* **6**(6), 861–867 (1993)
25. Mallat, S.: Group invariant scattering. *Commun. Pure Appl. Math.* **65**(10), 1331–1398 (2012)
26. Mallat, S.: Understanding deep convolutional networks. *Philos. Trans. R. Soc. A* **374**(2065), 20150203 (2016)
27. Manay, S., Cremers, D., Hong, B.-W., Yezzi, A.J., Soatto, S.: Integral invariants for shape matching. *IEEE Trans. Pattern Anal. Mach. Intell.* **28**(10), 1602–1618 (2006)
28. Marcos, D., Volpi, M., Komodakis, N., Tuia, D.: Rotation equivariant vector field networks. In: Proceedings of the IEEE International Conference on Computer Vision, pp. 5048–5057 (2017)
29. Mhaskar, H.N., Micchelli, C.A.: Approximation by superposition of sigmoidal and radial basis functions. *Adv. Appl. Math.* **13**(3), 350–373 (1992)
30. Munkres, J.R.: *Topology. Featured Titles for Topology Series*. Prentice Hall, Upper Saddle River (2000)
31. Pinkus, A.: TDI-subspaces of  $C(\mathbb{R}^d)$  and some density problems from neural networks. *J. Approx. Theory* **85**(3), 269–287 (1996)
32. Pinkus, A.: Approximation theory of the mlp model in neural networks. *Acta Numerica* **8**, 143–195 (1999)
33. Poggio, T., Mhaskar, H., Rosasco, L., Miranda, B., Liao, Q.: Why and when can deep-but not shallow-networks avoid the curse of dimensionality: a review. *Int. J. Autom. Comput.* 1–17 (2017)
34. Qi, C.R., Su, H., Mo, K., Guibas, L.J.: Pointnet: deep learning on point sets for 3d classification and segmentation. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pp. 652–660 (2017)
35. Reisert, M.: Group integration techniques in pattern analysis. Ph.D. thesis, Albert-Ludwigs-University (2008)
36. Schmid, B.J.: Finite groups and invariant theory. In: *Topics in Invariant Theory*, pp. 35–66. Springer (1991)
37. Schulz-Mirbach, H.: Invariant features for gray scale images. In: *Mustererkennung 1995*, pp. 1–14. Springer (1995)
38. Serre, J.-P.: *Linear Representations of Finite Groups*, vol. 42. Springer, Berlin (2012)
39. Sifre, L., Mallat, S.: Rigid-motion scattering for texture classification (2014). arXiv preprint [arXiv:1403.1687](https://arxiv.org/abs/1403.1687)
40. Simon, B.: *Representations of Finite and Compact Groups*. Number 10. American Mathematical Soc, London (1996)
41. Skibbe, H.: Spherical tensor algebra for biomedical image analysis. Ph.D. thesis, Albert-Ludwigs-University (2013)
42. Springenberg, J.T., Dosovitskiy, A., Brox, T., Riedmiller, M.: Striving for simplicity: the all convolutional net (2014). arXiv preprint [arXiv:1412.6806](https://arxiv.org/abs/1412.6806)
43. Thoma, M.: Analysis and optimization of convolutional neural network architectures. Masters’s thesis, Karlsruhe Institute of Technology, Karlsruhe, Germany, June 2017. <https://martin-thoma.com/msthesis/>

44. Vinberg, E.B.: *Linear Representations of Groups*. Birkhäuser, Basel (2012)
45. Waibel, A., Hanazawa, T., Hinton, G., Shikano, K., Lang, K.J.: Phoneme recognition using time-delay neural networks. *IEEE Trans. Acoust. Speech Signal Process.* **37**(3), 328–339 (1989)
46. Weyl, H.: *The Classical Groups: Their Invariants and Representations*. Princeton Mathematical Series (1) (1946)
47. Worfolk, P.A.: Zeros of equivariant vector fields: algorithms for an invariant approach. *J. Symb. Comput.* **17**(6), 487–511 (1994)
48. Worrall, D.E., Garbin, S.J., Turmukhambetov, D., 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)
49. Zaheer, M., Kottur, S., Ravanbakhsh, S., Poczos, B., Salakhutdinov, R.R., Smola, A.J.: Deep sets. In: *Advances in Neural Information Processing Systems*, pp. 3391–3401 (2017)

**Publisher's Note** Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.