# UNIFIED UNIVERSALITY THEOREM FOR DEEP AND SHALLOW JOINT-GROUP-EQUIVARIANT MACHINES

Anonymous authors

004

010 011

012

013

014

015

016

017

018

019

021

023

040

Paper under double-blind review

#### ABSTRACT

We present a constructive universal approximation theorem for learning machines equipped with joint-group-equivariant feature maps, based on the group representation theory. "Constructive" here indicates that the distribution of parameters is given in a closed-form expression known as the ridgelet transform. Joint-groupequivariance encompasses a broad class of feature maps that generalize classical group-equivariance. Notably, this class includes fully-connected networks, which are *not* group-equivariant *but* are joint-group-equivariant. Moreover, our main theorem also unifies the universal approximation theorems for both shallow and deep networks. While the universality of shallow networks has been investigated in a unified manner by the ridgelet transform, the universality of deep networks has been investigated in a case-by-case manner.

#### 1 INTRODUCTION

025 The proof of a universality theorem contains hints for understanding the internal data processing 026 mechanisms inside learning machines such as neural networks. For example, the first universality theorems for depth-2 neural networks were shown in 1989 with four different proofs by Cybenko 027 028 (1989), Hornik et al. (1989), Funahashi (1989), and Carroll & Dickinson (1989). Among them, Cy-029 benko's proof using Hahn-Banach and Hornik et al.'s proof using Stone-Weierstrass are existential proofs, meaning that it is not clear how to assign the parameters. On the other hand, Funahashi's proof reducing to the Fourier transform and Carroll and Dickinson's proof reducing to the Radon 031 transform are constructive proofs, meaning that it is clear how to assign the parameters. The lat-032 ter constructive methods, which reduce to integral transforms, were refined as the so-called integral 033 representation by Barron (1993) and further culminated as the *ridgelet transform*, the main objective 034 of this study, discovered by Murata (1996) and Candès (1998).

To show the universality in a constructive manner, we formulate the problem as a functional equation: Let  $LM[\gamma]$  denote a certain learning machine (such as a deep network) with parameter  $\gamma$ , and let  $\mathcal{F}$  denote a class of functions to be expressed by the learning machine. Given a function  $f \in \mathcal{F}$ , find an unknown parameter  $\gamma$  so that the machine  $LM[\gamma]$  represents function f, i.e.

 $LM[\gamma] = f, \tag{1}$ 

041 which we call a *learning equation*. This equation is understood as a stronger formulation of 042 learning than an ordinary formulation by the empirical risk minimization such as minimizing  $\sum_{i=1}^{n} |\mathrm{LM}[\gamma](x_i) - f(x_i)|^2$  with respect to  $\gamma$ , as the latter is understood as a weak form (or a varia-043 044 tional form) of this equation. Therefore, characterizing the solution space of this equation leads to understanding the parameters obtained by risk minimization. Following previous studies (Murata, 1996; Candès, 1998; Sonoda et al., 2021a;b; 2022a;b), we call a solution operator R that satisfies 046 LM[R[f]] = f a ridgelet transform. Once such a solution operator R is found, we can conclude a uni-047 *versality* of the learning machine in consideration because the reconstruction formula LM[R[f]] = f048 implies for any  $f \in \mathcal{F}$  there exists a machine that represents f. In particular, when R[f] is found in a closed-form manner, then it leads to a *constructive* proof of the universality since R[f] could indicate how to assign parameters. 051

For depth-2 neural networks (particularly with an infinitely-wide hidden layer), the equation has
 been solved with several closed-form ridgelet transforms. For example, the closed-form ridgelet transforms have been obtained for depth-2 fully-connected layers (Sonoda et al., 2021b), depth-2

fully-connected layers on manifolds (Sonoda et al., 2022b), depth-2 group convolution layers (Sonoda et al., 2022a), and depth-2 fully-connected layers on finite fields (Yamasaki et al., 2023). The essential technique to obtain these ridgelet transforms are to construct a Fourier expression corresponding to the network in consideration. We refer to Sonoda et al. (2024b) for more tecnnical backgrounds behind these results. Furthermore, Sonoda et al. (2021a) have revealed that the distribution of parameters inside depth-2 fully-connected networks obtained by regularized empirical risk minimization assymptotically converges to the ridgelet transform. In other words, the ridgelet transform can also explain the solutions obtained by risk minimization.

062 On the other hand, for depth-n neural networks, the equation is far from solved, and it is common 063 to either consider infinitely-deep mathematical models such as Neural ODEs (Sonoda & Murata, 064 2017b; E, 2017; Li & Hao, 2018; Haber & Ruthotto, 2017; Chen et al., 2018), or handcraft solutions depending on the network specifications. For example, construction methods such as the so-called 065 Telgarsky sawtooth function (or the Yarotsky scheme) and bit extraction techniques (Cohen et al., 066 2016; Telgarsky, 2016; Yarotsky, 2017; 2018; Yarotsky & Zhevnerchuk, 2020; Daubechies et al., 067 2022; Cohen et al., 2022; Siegel, 2023; Petrova & Wojtaszczyk, 2023; Grohs et al., 2023) have been 068 developed (not only to investigate the expressivity but also) to demonstrate the depth separation, 069 super-convergence, and minmax optimality of deep ReLU networks. Various feature maps have also been handcrafted in the contexts of geometric deep learning (Bronstein et al., 2021) and deep narrow 071 networks (Lu et al., 2017; Hanin & Sellke, 2017; Lin & Jegelka, 2018; Kidger & Lyons, 2020; Park et al., 2021; Li et al., 2023; Cai, 2023; Kim et al., 2024). However, for the purpose of understanding 073 the parameters obtained by risk minimization (in a manner presented by Sonoda et al. (2021a)), 074 these results are less satisfactory because there is no guarantee that these handcrafted solutions are 075 obtaiend by risk minimization.

076 Recently, Sonoda et al. (2024a) developed a novel technique to show the universality based on the 077 group representation theory, and discovered a rich class of ridgelet transforms for learning machines 078 with joint-group-invariant feature maps. However, their technique was essentially limited to depth-079 2 networks and could not cover depth-n networks defined by composites of nonlinear activation 080 functions such as  $\sigma(A_2\sigma(A_1x - b_1) - b_2)$ . By carefully reviewing their group theoretic arguments, 081 we found that the joint-invariance is the bottleneck and it can be resolved by relaxing the assumption 082 to the joint-equivariance. In this study, we present a wider class of ridgelet transforms for learning 083 machines with joint-equivariant feature maps so to cover the depth-n (as well as depth-2) fullyconnected networks. 084

- 085 The contributions of this study include
  - We derived the ridgelet transform (solution operator for the learning equation) for learning machines with depth-n joint-group-equivariant feature maps. Since the solution of the learning equation can be written in closed form for any  $f \in \mathcal{F}$ , it is a constructive and unified proof of the universal approximation theorem for joint-group-equivariant machines.
  - As a corollary, we have shown the constructive universal approximation property of deep fully-connected neural networks. Until this study, the universality of deep networks has been shown in a different manner from the universality of shallow networks, but our results discuss them on common ground. Now we can understand the approximation schemes of various learning machines in a unified manner.
    - In addition, as an example of a learning machine whose universality has not been known, we presented a network with quadratic forms and showed its universality.
    - Further, we have shown the ridgelet transform for depth-*n* group convolutional networks. In a previous study, it was known only for depth-2, and we have succeeded to extend it by reviewing the arguments from the group theoretic perspective.
- 100 101 102

087

090

092

093

094

095 096

098

099

### 2 PRELIMINARIES

103 104

We quickly overview the original integral representation and the ridgelet transform, a mathematical
 model of depth-2 fully-connected network and its right inverse. Then, we list a few facts in the
 group representation theory. In particular, *Schur's lemma* and the *Haar measure* play key roles in
 the proof of the main results.

**Notation.** For any topological space X,  $C_c(X)$  denotes the Banach space of all compactly supported continuous functions on X. For any measure space X,  $L^p(X)$  denotes the Banach space of all *p*-integrable functions on X.  $S(\mathbb{R}^d)$  and  $S'(\mathbb{R}^d)$  denote the classes of rapidly decreasing functions (or Schwartz test functions) and tempered distributions on  $\mathbb{R}^d$ , respectively.

#### 2.1 INTEGRAL REPRESENTATION AND RIDGELET TRANSFORM FOR DEPTH-2 FULLY-CONNECTED NETWORK

**Definition 1.** For any measurable functions  $\sigma : \mathbb{R} \to \mathbb{C}$  and  $\gamma : \mathbb{R}^m \times \mathbb{R} \to \mathbb{C}$ , put

$$S_{\sigma}[\gamma](\boldsymbol{x}) := \int_{\mathbb{R}^m \times \mathbb{R}} \gamma(\boldsymbol{a}, b) \sigma(\boldsymbol{a} \cdot \boldsymbol{x} - b) \mathrm{d}\boldsymbol{a} \mathrm{d}b, \quad \boldsymbol{x} \in \mathbb{R}^m.$$
(2)

We call  $S_{\sigma}[\gamma]$  an (integral representation of) neural network, and  $\gamma$  a parameter distribution.

The integration over all the hidden parameters  $(a, b) \in \mathbb{R}^m \times \mathbb{R}$  means all the neurons  $\{x \mapsto \sigma(a \cdot x - b) \mid (a, b) \in \mathbb{R}^m \times \mathbb{R}\}$  are summed (or integrated, to be precise) with weight  $\gamma$ , hence formally  $S_{\sigma}[\gamma]$  is understood as a continuous neural network with a single hidden layer. We note, however, when  $\gamma$  is a finite sum of point measures such as  $\gamma_p = \sum_{i=1}^p c_i \delta_{(a_i,b_i)}$  (by appropriately extending the class of  $\gamma$  to Borel measures), then it can also reproduce a finite width network

$$S_{\sigma}[\gamma_p](\boldsymbol{x}) = \sum_{i=1}^p c_i \sigma(\boldsymbol{a}_i \cdot \boldsymbol{x} - b_i).$$
(3)

In other words, the integral representation is a mathematical model of depth-2 network with *any* width (ranging from finite to continuous).

132 Next, we introduce the ridgelet transform, which is known to be a right-inverse operator to  $S_{\sigma}$ .

**Definition 2.** For any measurable functions  $\rho : \mathbb{R} \to \mathbb{C}$  and  $f : \mathbb{R}^m \to \mathbb{C}$ , put

$$R_{\rho}[f](\boldsymbol{a},b) := \int_{\mathbb{R}^m} f(\boldsymbol{x}) \overline{\rho(\boldsymbol{a} \cdot \boldsymbol{x} - b)} \mathrm{d}\boldsymbol{x}, \quad (\boldsymbol{a},b) \in \mathbb{R}^m \times \mathbb{R}.$$
(4)

137 We call  $R_{\rho}$  a ridgelet transform.

To be precise, it satisfies the following reconstruction formula.

**Theorem 1** (Reconstruction Formula). Suppose  $\sigma$  and  $\rho$  are a tempered distribution (S') and a rapid decreasing function (S) respectively. There exists a bilinear form  $((\sigma, \rho))$  such that

 $S_{\sigma} \circ R_{\rho}[f] = ((\sigma, \rho))f, \tag{5}$ 

for any square integrable function  $f \in L^2(\mathbb{R}^m)$ . Further, the bilinear form is given by  $((\sigma, \rho)) = \int_{\mathbb{R}} \sigma^{\sharp}(\omega) \overline{\rho^{\sharp}(\omega)} |\omega|^{-m} d\omega$ , where  $\sharp$  denotes the 1-dimensional Fourier transform.

146 See Sonoda et al. (2021b, Theorem 6) for the proof. In particular, according to Sonoda et al. (2021b, 147 Lemma 9), for any activation function  $\sigma$ , there always exists  $\rho$  satisfying  $((\sigma, \rho)) = 1$ . Here,  $\sigma$ 148 being a tempered distribution means that typical activation functions are covered such as ReLU, 149 step function, tanh, gaussian, etc... We can interpret the reconstruction formula as a universality 150 theorem of continuous neural networks, since for any given data generating function f, a network 151 with output weight  $\gamma_f = R_{\rho}[f]$  reproduces f (up to factor  $((\sigma, \rho))$ ), i.e.  $S[\gamma_f] = f$ . In other words, the ridgelet transform indicates how the network parameters should be organized so that the network 152 represents an individual function f. 153

The original ridgelet transform was discovered by Murata (1996) and Candès (1998). It is recently extended to a few modern networks by the Fourier slice method (see e.g. Sonoda et al., 2024b).
In this study, we present a systematic scheme to find the ridgelet transform for a variety of given network architecture based on the group theoretic arguments.

158 159

160

112 113

114

115

120

127 128

141 142

143

2.2 IRREDUCIBLE UNITARY REPRESENTATION AND SCHUR'S LEMMA

In the main theorem, we use *Schur's lemma*, a fundamental theorem from unitary group representation theory. Group representation is a method for investigating properties of an abstract group G by mapping G to another (much computable) group of invertible linear operators. We refer to Folland (2015) for more details on group representation and harmonic analysis on groups.

In this study, we assume group G to be *locally compact*. This is a sufficient condition for having invariant measures. It is not a strong assumption. For example, any finite group, discrete group, compact group, and finite-dimensional Lie group are locally compact, while an infinite-dimensional Lie group is *not* locally compact.

Let  $\mathcal{H}$  be a nonzero Hilbert space, and  $\mathcal{U}(\mathcal{H})$  be the group of unitary operators on  $\mathcal{H}$ . A *unitary* representation  $\pi$  of G on  $\mathcal{H}$  is a group homomorphism that is continuous with respect to the strong operator topology—that is, a map  $\pi : G \to \mathcal{U}(\mathcal{H})$  satisfying  $\pi_{gh} = \pi_g \pi_h$  and  $\pi_{g^{-1}} = \pi_g^{-1}$ , and for any  $\psi \in \mathcal{H}$ , the map  $G \ni g \mapsto \pi_g[\psi] \in \mathcal{H}$  is continuous.

173 Suppose  $\mathcal{M}$  is a closed subspace of  $\mathcal{H}$ .  $\mathcal{M}$  is called an *invariant* subspace when  $\pi_g[\mathcal{M}] \subset \mathcal{M}$  for all 174  $g \in G$ . Particularly,  $\pi$  is called *irreducible* when it does not admit any nontrivial invariant subspace 175  $\mathcal{M} \neq \{0\}$  nor  $\mathcal{H}$ . The following theorem is a fundamental result of group representation theory that 176 characterizes the irreducibility.

**Theorem 2** (Schur's lemma). A unitary representation  $\pi : G \to U(\mathcal{H})$  is irreducible iff any bounded operator T on  $\mathcal{H}$  that commutes with  $\pi$  is always a constant multiple of the identity. In other words, if  $\pi_g \circ T = T \circ \pi_g$  for all  $g \in G$ , then  $T = c \operatorname{Id}_{\mathcal{H}}$  for some  $c \in \mathbb{C}$ .

See Folland (2015, Theorem 3.5(a)) for the proof. We use this as a key step in the proof of our main theorem.

As a concrete example of an irreducible representation, we use the following regular representation of the affine group  $\operatorname{Aff}(m)$  on  $L^2(\mathbb{R}^m)$ .

**Theorem 3.** Let  $G := \operatorname{Aff}(m) := GL(m) \ltimes \mathbb{R}^m$  be the affine group acting on  $X = \mathbb{R}^m$  by (L,t)  $\cdot x = Lx + t$ , and let  $\mathcal{H} := L^2(\mathbb{R}^m)$  be the Hilbert space of square-integrable functions. Let  $\pi : \operatorname{Aff}(m) \to \mathcal{U}(L^2(\mathbb{R}^m))$  be the regular representation of the affine group  $\operatorname{Aff}(m)$  on  $L^2(\mathbb{R}^m)$ , namely  $\pi_q[f](x) := |\det L|^{-1/2} f(L^{-1}(x-t))$  for any  $g = (L,t) \in G$ . Then  $\pi$  is irreducible.

189 190 See Folland (2015, Theorem 6.42) for the proof.

#### 2.3 CALCULUS ON LOCALLY COMPACT GROUP

By Haar's theorem, if G is a locally compact group, then there uniquely exist left and right invariant measures  $d_l g$  and  $d_r g$ , satisfying for any  $s \in G$  and  $f \in C_c(G)$ ,

$$\int_{G} f(sg) \mathrm{d}_{l}g = \int_{G} f(g) \mathrm{d}_{l}g, \quad \text{and} \quad \int_{G} f(gs) \mathrm{d}_{r}g = \int_{G} f(g) \mathrm{d}_{r}g.$$

Let X be a G-space with transitive left (resp. right) G-action  $g \cdot x$  (resp.  $x \cdot g$ ) for any  $(g, x) \in G \times X$ . Then, we can further induce the left (resp. right) invariant measure  $d_l x$  (resp.  $d_r x$ ) so that for any  $f \in C_c(G)$ ,

$$\int_X f(x) \mathrm{d}_l x := \int_G f(g \cdot o) \mathrm{d}_l g, \quad \text{resp.} \quad \int_X f(x) \mathrm{d}_r x := \int_G f(o \cdot g) \mathrm{d}_r g,$$

where  $o \in X$  is a fixed point called the origin.

#### 3 MAIN RESULTS

191

192 193

194

195 196 197

199

200

205 206

207 208

We introduce unitary representations  $\pi$  and  $\hat{\pi}$ , a *joint-equivariant feature map*  $\phi : X \times \Xi \to Y$ , a *joint-equivariant machine*  $LM[\gamma; \phi] : X \to Y$ , and present the ridgelet transform  $R[f; \psi] : \Xi \to \mathbb{C}$ for joint-equivariant machines, yielding the universality  $LM[R[f; \psi]; \phi] = c_{\phi,\psi}f$ . We note that  $\pi$ plays a key role in the main theorem, and the joint-equivariance is an essential property of depth-*n* fully-connected network.

Let G be a locally compact group equipped with a left invariant measure dg. Let X and  $\Xi$  be Gspaces equipped with G-invariant measures dx and d $\xi$ , called the *data domain* and the *parameter* domain. respectively. Let Y be a separable Hilbert space, called the *output domain*. Let  $\mathcal{U}(Y)$  be



Figure 1: The classical G-equivariant feature map  $\phi: X \times \Xi \rightarrow$  Figure 2: Joint-equivariant map  $\phi$ Y is a subclass of joint-G-equivariant map where the G-action is a G-equivariant section of Gbundle over base  $\Xi$  with fiber  $Y^X$ on parameter domain  $\Xi$  is *trivial*, i.e.  $g \cdot \xi = \xi$ 

the space of unitary operators on Y, and let  $v: G \to \mathcal{U}(Y)$  be a unitary representation of G on Y. We call a Y-valued map  $\phi$  on the data-parameter domain  $X \times \Xi$ , i.e.  $\phi : X \times \Xi \to Y$ , a *feature* map.

232 Let  $L^2(X;Y)$  denote the space of Y-valued square-integrable functions on X equipped with the 233 inner product  $\langle \phi, \psi \rangle_{L^2(X;Y)} := \int_X \langle \phi(x), \psi(x) \rangle_Y dx$ ; and let  $L^2(\Xi)$  denote the space of  $\mathbb{C}$ -valued 234 square-integrable functions on  $\Xi$ . 235

If there is no risk of confusion, we use the same symbol  $\cdot$  for the G-actions on X, Y, and  $\Xi$  (e.g., 236  $g \cdot x, g \cdot y$ , and  $g \cdot \xi$ ). On the other hand, to avoid the confusion between G-actions on output domain 237 Y and Y-valued function  $f: X \to Y$ , both " $g \cdot f(x)$ " and " $v_q[f(x)]$ " (if needed) always imply 238 G-action on Y, and " $\pi_q[f](x)$ " (introduced soon below) for G-actions on  $f: X \to Y$ . 239

Additionally, we introduce two unitary representations  $\pi$  and  $\hat{\pi}$  of G on function spaces  $L^2(X;Y)$ 240 and  $L^2(\Xi)$  as follows: For each  $g \in G$ ,  $f \in L^2(X;Y)$  and  $\gamma \in L^2(\Xi)$ , 241

$$\pi_g[f](x) := v_g[f(g^{-1} \cdot x)] = g \cdot f(g^{-1} \cdot x), \quad x \in X$$
(6)

247

248 249

250

256 257

225

226

227 228 229

230

231

$$\widehat{\pi}_q[\gamma](\xi) := \gamma(g^{-1} \cdot \xi), \quad \xi \in \Xi.$$
(7)

In the main theorem, the irreducibility of  $\pi$  will be a sufficient condition for the universality. On 246 the other hand, the irreducibility of  $\hat{\pi}$  is not necessary. For those who are less familiar with group representations, we have shown that  $\pi$  and  $\hat{\pi}$  are unitary representations in Lemmas 6 and 7.

#### 3.1 JOINT-EQUIVARIANT FEATURE MAP

251 We introduce the joint-group-equivariant feature map, extending the classical notion of group-252 equivariant feature maps. The major motivation to introduce this is that depth-n fully-connected networks, the main subject of this study, are not equivariant but joint-equivariant. 253

254 **Definition 3** (Joint-G-Equivariant Feature Map). We say a feature map  $\phi: X \times \Xi \to Y$  is joint-G-255 *equivariant* when

$$\phi(g \cdot x, g \cdot \xi) = g \cdot \phi(x, \xi), \quad (x, \xi) \in X \times \Xi,$$
(8)

258 holds for all  $g \in G$ . Especially, when G-action on Y is trivial, i.e.  $\phi(g \cdot x, g \cdot \xi) = \phi(x, \xi)$ , we say it is *joint-G-invariant*. 259

260 Remark 1 (Relation to classical G-equivariance). The joint-G-equivariance is not a restriction but 261 an extension of the classical notion of G-equivariance, i.e.  $\phi(g \cdot x, \xi) = g \cdot \phi(x, \xi)$ . In fact, G-262 equivariance is a special case of joint-G-equivariance where G acts trivially on parameter domain, 263 i.e.  $g \cdot \xi = \xi$  (see Figure 1). Thus, all G-equivariant maps are automatically joint-G-equivariant.

264 *Remark* 2 (Interpretations of joint-G-equivariance). We have two interpretations from algebraic and 265 geometric perspectives. First, from an algebraic perspective,  $\phi$  is a homomorphism (or a G-map) 266 between G-sets from  $X \times \Xi$  to Y. So we may denote the collection of all joint-G-equivariant maps as  $\hom_G(X \times \Xi, Y)$ . Second, from a more geometric perspective,  $\phi$  is a vector-field  $\Xi \to Y^X$  with 267 structure group G acting on fiber  $Y^X$  by  $\pi$ . Here we identify  $\phi: X \times \Xi \to Y$  with  $\phi_c: \Xi \to Y^X$ 268 by the so-called *currying*  $\phi_c(\xi) := \phi(\bullet, \xi)$ . In other words,  $\phi_c$  is a global section of a trivial G-269 bundle  $\Xi \times Y^X \to \Xi$ . Consequently, we can understand the G-action  $g \cdot \xi$  on parameter domain (or

base space)  $\Xi$  is induced from the coordinate transformation  $\pi$  on fiber  $Y^X$  so that the section  $\phi_c$  is *G*-equivariant (see Figure 2)

$$\pi_{q}[\phi_{c}(\xi)](x) = g \cdot \phi(g^{-1} \cdot x, \xi) =: \phi_{c}(g \cdot \xi)(x).$$
(9)

Finally, two aspects are unified as a tensor-hom adjunction:  $\hom_G(X \times \Xi, Y) \cong \hom_G(\Xi, Y^X)$ .

In the following, we list several construction methods of joint-equivariant maps in Lemmas 1, 2 and (in the next subsection), indicating the richness of the proposed concept. Whereas to construct a (non-joint) *G*-equivariant network, we must carefully and precisely design the network architecture (see, e.g., a textbook of geometric deep learning Bronstein et al., 2021), to construct a joint-*G*equivariant network, we can easily and systematically obtain the one.

First, we can synthesize a joint-equivariant map from (not equivariant but) any map  $\phi_0: X \to Y$ .

**Lemma 1.** Let X and Y be G-sets. Fix an arbitrary map  $\phi_0 : X \to Y$ , and put  $\phi(x,g) := \pi_g[\phi_0](x) = g \cdot \phi_0(g^{-1} \cdot x)$  for every  $x \in X$  and  $g \in G$ . Then,  $\phi : X \times G \to Y$  is joint-G-equivariant.

*Proof.* For any 
$$g, h \in G$$
, we have  $\phi(g \cdot x, g \cdot h) = (gh) \cdot \phi_0((gh)^{-1} \cdot (g \cdot x)) = g \cdot \phi(x, h)$ .  $\Box$ 

In general, a G-set is understood as a representation of G. So, the case of  $X = Y = \Xi = G$  with  $\phi: G \times G \to G$  is understood as a primitive type of joint-G-equivariant maps  $\phi: X \times \Xi \to Y$ .

291 The next lemma suggests the compatibility with function compositions, or deep structures.

**Lemma 2** (Depth-*n* Joint-Equivariant Feature Map  $\phi_{1:n}$ ). *Given a sequence of joint-G-equivariant* feature maps  $\phi_i : X_{i-1} \times \Xi_i \to X_i$  (i = 1, ..., n), let  $\Xi_{1:n} := \Xi_1 \times \cdots \times \Xi_n$  be the *n*-fold parameter space with the component-wise *G*-action  $g \cdot \xi_{1:n} := (g \cdot \xi_1, ..., g \cdot \xi_n)$  for each *n*-fold parameters  $\xi_{1:n} \in \Xi_{1:n}$ , and let  $\phi_{1:n} : X_0 \times \Xi_{1:n} \to X_n$  be the depth-*n* feature map given by

$$\phi_{1:n}(x,\xi_{1:n}) := \phi_n(\bullet,\xi_n) \circ \dots \circ \phi_1(x,\xi_1).$$

$$(10)$$

298 Then,  $\phi_{1:n}$  is joint-G-equivariant. 299

> In other words, the composition of joint-equivariant maps defines a cascade product of morphisms:  $\hom_G(\Xi_2, X_2^{X_1}) \times \hom_G(\Xi_1, X_1^{X_0}) \to \hom_G(\Xi_1 \times \Xi_2, X_2^{X_0})$ . See Appendix A.2 for the proof.

#### 3.2 JOINT-EQUIVARIANT MACHINE AND RIDGELET TRANSFORM

We further introduce the joint-equivariant machine, extending the integral representation.

**Definition 4** (Joint-Equivariant Machine). Fix an arbitrary joint-equivariant feature map  $\phi : X \times \Xi \to Y$ . For any scalar-valued measurable function  $\gamma : \Xi \to \mathbb{C}$ , define a Y-valued map on X by

$$LM[\gamma;\phi](x) := \int_{\Xi} \gamma(\xi)\phi(x,\xi)d\xi, \quad x \in X,$$
(11)

where the integral is understood as the Bochner integral. We also write  $LM_{\phi} := LM[\bullet; \phi]$  for short. If needed, we call the image  $LM[\gamma; \phi] : X \to Y$  a joint-equivariant *machine*, and the integral transform  $LM[\bullet; \phi]$  of  $\gamma$  a joint-equivariant *transform*.

The joint-equivariant machine extends the original integral representation. It inherits the concept of integrating all the possible parameters  $\xi$  and indirectly select which parameters to use by weighting on them, which *linearize* parametrization by lifting nonlinear parameters  $\xi$  to linear parameter  $\gamma$ .

Recall that the *G*-action on parameter domain  $\Xi$  is also linearized by lifting it to  $\hat{\pi}$  on  $L^2(\Xi)$ . The joint-equivariance of  $\phi : \Xi \to Y^X$  is inherited under the linearization to  $LM_{\phi} : L^2(\Xi) \to L^2(X;Y)$ .

Lemma 3. A joint-G-equivariant machine  $LM_{\phi} : L^2(\Xi) \to L^2(X;Y)$  is joint-G-equivariant, i.e.  $LM_{\phi} \in \hom_G(L^2(\Xi), L^2(X;Y)).$ 

323

273

286 287 288

296 297

300

301 302 303

304

305 306

307

308 309 310

314

$$Proof. \ \mathsf{LM}_{\phi}[\widehat{\pi}_{g}[\gamma]](g \cdot x) = \int_{\Xi} \gamma(g^{-1} \cdot \xi) \phi(g \cdot x, \xi) \mathrm{d}\xi = \int_{\Xi} \gamma(\xi) \phi(g \cdot x, g \cdot \xi) \mathrm{d}\xi = g \cdot \mathsf{LM}_{\phi}[\gamma](x). \quad \Box$$

**Definition 5** (Ridgelet Transform for Joint-Equivariant Machine). For any joint-equivariant feature map  $\psi : X \times \Xi \to Y$  and Y-valued Borel measurable function f on X, put a scalar-valued map by

$$\mathbb{R}[f;\psi](\xi) := \int_X \langle f(x), \psi(x,\xi) \rangle_Y \mathrm{d}x, \quad \xi \in \Xi.$$
(12)

 We also write  $\mathbb{R}_{\psi} := \mathbb{R}[\bullet; \psi]$  for short. If there is no risk of confusion, we call both the image  $\mathbb{R}[f; \psi] : X \to Y$  and the integral transform  $\mathbb{R}[\bullet; \psi]$  of f a ridgelet transform.

Intuitively, it measures the similarity between target function f and feature  $\psi(\bullet, \xi)$  at  $\xi$ . As long as the integrals are convergent, the ridgelet transform is the dual operator of the joint-equivariant transform (with common  $\phi$ ):

$$\langle \gamma, \mathbb{R}[f;\phi] \rangle_{L^2(\Xi)} = \int_{X \times \Xi} \gamma(\xi) \langle \phi(x,\xi), f(x) \rangle_Y \mathrm{d}x \mathrm{d}\xi = \langle \mathbb{L}\mathbb{M}[\gamma;\phi], f \rangle_{L^2(X;Y)}.$$
(13)

Similarly to the joint-equivariant machine, the ridgelet transform is again joint-G-invariant. In fact,

$$\mathbf{R}_{\psi}[\pi_{g}[f]](g \cdot \xi) = \int_{X} \langle v_{g}[f(g^{-1} \cdot x)], \psi(x, g \cdot \xi) \rangle_{Y} \mathrm{d}\xi = \int_{X} \langle f(x), v_{g}^{*}[\psi(g \cdot x, g \cdot \xi)] \rangle_{Y} \mathrm{d}\xi = \mathbf{R}_{\psi}[f](\xi)$$

Hence, geometrically, if we regard  $LM_{\phi} : L^2(\Xi) \to L^2(X;Y)$  a vector field of trivial *G*-bundle  $L^2(\Xi) \times L^2(X;Y) \to L^2(\Xi)$ , then  $R_{\phi} : L^2(X;Y) \to L^2(\Xi)$  corresponds to a *G*-connection.

#### 3.3 MAIN THEOREM

At last, we state the main theorem, that is, the reconstruction formula for joint-equivariant machines. Theorem 4 (Reconstruction Formula). Assume (1) feature maps  $\phi, \psi : X \times \Xi \rightarrow Y$  are joint-Gequivariant, (2) composite operator  $LM_{\phi} \circ R_{\psi} : L^2(X;Y) \rightarrow L^2(X;Y)$  is bounded (i.e., Lipschitz continuous), and (3) the unitary representation  $\pi : G \rightarrow U(L^2(X;Y))$  defined in (6) is irreducible. Then, there exists a bilinear form  $((\phi, \psi)) \in \mathbb{C}$  (independent of f) such that for any Y-valued squareintegrable function  $f \in L^2(X;Y)$ ,

$$\mathrm{LM}_{\phi} \circ \mathbf{R}_{\psi}[f] = \int_{\Xi} \left[ \int_{X} \langle f(x), \psi(x,\xi) \rangle_{Y} \mathrm{d}x \right] \phi(\bullet,\xi) \mathrm{d}\xi = ((\phi,\psi))f.$$
(14)

In practice, once the irreducibility of *G*-action  $\pi$  on  $L^2(X;Y)$  is verified, the ridgelet transform R<sub>\u03c0</sub> becomes a right inverse operator of joint-equivariant transform LM<sub>\u03c0</sub> as long as  $((\u03c0, \u03c0)) \neq 0, \infty$ . Despite the wide coverage of examples, the proof is brief and simple as follows.

*Proof.* By using the unitarity of representation  $v : G \to U(Y)$ , left-invariance of measure dx, and 361 *G*-equivariance of feature map  $\psi$ , for all  $g \in G$ , we have

$$\mathbf{R}_{\psi}[\pi_{g}[f]](\xi) = \int_{X} \langle g \cdot f(g^{-1} \cdot x), \psi(x,\xi) \rangle_{Y} \mathrm{d}x = \int_{X} \langle f(x), g^{-1} \cdot \psi(g \cdot x,\xi) \rangle_{Y} \mathrm{d}x \\
= \int_{X} \langle f(x), \psi(x, g^{-1} \cdot \xi) \rangle_{Y} \mathrm{d}x = \widehat{\pi}_{g}[\mathbf{R}_{\psi}[f]](\xi).$$
(15)

Similarly,

$$LM_{\phi}[\widehat{\pi}_{g}[\gamma]](x) = \int_{\Xi} \gamma(g^{-1} \cdot \xi)\phi(x,\xi)d\xi = \int_{\Xi} \gamma(\xi)\phi(x,g \cdot \xi)d\xi$$
$$= \int_{\Xi} \gamma(\xi) \left(g \cdot \phi(g^{-1} \cdot x,\xi)\right)d\xi = \pi_{g}[LM_{\phi}[\gamma]](x).$$
(16)

As a consequence,  $LM_{\phi} \circ R_{\psi} : L^2(X;Y) \to L^2(X;Y)$  commutes with  $\pi$  as below

$$\mathsf{LM}_{\phi} \circ \mathsf{R}_{\psi} \circ \pi_{g} = \mathsf{LM}_{\phi} \circ \widehat{\pi}_{g} \circ \mathsf{R}_{\psi} = \pi_{g} \circ \mathsf{LM}_{\phi} \circ \mathsf{R}_{\psi}$$
(17)

for all  $g \in G$ . Hence by Schur's lemma (Theorem 2), there exist a constant  $C_{\phi,\psi} \in \mathbb{C}$  such that  $LM_{\phi} \circ R_{\psi} = C_{\phi,\psi} Id_{L^2(X)}$ . Since  $LM_{\phi} \circ R_{\psi}$  is bilinear in  $\phi$  and  $\psi$ ,  $C_{\phi,\psi}$  is bilinear in  $\phi$  and  $\psi$ .  $\Box$ 



Figure 3: Deep Y-valued joint-G-equivariant machine on G-space X is  $L^2(X; Y)$ -universal when unitary representation  $\pi$  of G on  $L^2(X; Y)$  is irreducible, and the distribution of parameters for the machine to represent a given map  $f: X \to Y$  is exactly given by the ridgelet transform  $\mathbb{R}[f]$ 

391 *Remark* 3. (1) When  $\pi$  is not irreducible (thus reducible) and admits an irreducible decomposition 392  $L^2(X;Y) = \bigoplus_{i=1}^{\infty} \mathcal{H}_i$ , reconstruction formula (14) holds for  $f \in \mathcal{H}_k$  for some k. This is another 393 consequence from Schur's lemma. (2) The irreducibility is assumed only for  $\pi$ , and not for  $\hat{\pi}$ . This asymmetry originates from the fact that our main theorem focuses only on (the universality of) 394 395  $\mathrm{LM}_{\phi}[\gamma]: X \to Y$ , not on its dual  $\mathrm{R}_{\psi}[f]: \Xi \to \mathbb{R}$ . However, when  $\widehat{\pi}$  is irreducible, then we can state  $\mathbb{R}_{\psi} \circ \mathbb{L}\mathbb{M}_{\phi}[\gamma] = \gamma$  for any  $\gamma \in L^2(\Xi)$  (the order of composition is reverted from  $\mathbb{L}\mathbb{M}_{\phi} \circ \mathbb{R}_{\psi}$ ). (3) The regularity of feature maps  $\phi, \psi$  needs to be studied in a case-by-case manner. Showing the 397 joint-equivariance is relatively easy. For example, for fully-connected networks (§ 5) and quadratic-398 form networks ( $\S$  6), the joint-equivariance holds for any activation function. However, the constant 399  $((\phi, \psi))$  can degenerate to zero or diverge depending on feature maps. For the case of depth-2 fully-400 connected networks, it is known that the constant is zero if and only if the activation function is 401 a polynomial function (see e.g., Sonoda & Murata, 2017a). In general, such a condition can be 402 investigated in a case-by-case manner. Fortunately, we can use the closed-form expression of the 403 ridgelet transform to our advantage.

404 Remark 4 (Technical differences from Sonoda et al. (2024a, Theorem 9)). The previous study cannot 405 deal with deep structures or composite maps and, therefore, falls short as a theoretical analysis 406 of deep learning. This limitation arises because joint-invariance alone is insufficient to construct 407 irreducible representations (or irreps, for short) for deep structures. More technically, it is due to 408 the facts (1) that an *inner* tensor product  $\pi_1 \otimes_i \pi_2$  of irreps  $\pi_1, \pi_2$  is not always irreducible (see e.g. Chapter 7.3 Folland, 2015), while (2, or Lemma 4) that an *outer* tensor product  $\pi_1 \otimes \pi_2$  or 409 irreps  $\pi_1, \pi_2$  is always irreducible. The two facts (1) and (2) appear similar (thus confusing) but the 410 essential consequences are different. Further, (3) that deep feature maps are often vector-valued, but 411 (4) that the previous study is limited to scalar-valued joint-invariant feature maps. In the previous 412 study, due to (1) and (4), it is technically hard to obtain an irrep for vector-valued feature maps. 413 Namely, just a d-times inner tensor product  $\pi_s \otimes_i \cdots \otimes_i \pi_s$  of irrep  $\pi_s$  acting on scalar-valued joint-414 invariant maps  $\phi_s$  cannot be an irrep acting on d-dim. vector-valued feature map  $\phi_s \otimes \cdots \otimes \phi_s$ . On 415 the other hand, this study is based on not (1) but (2), we have successfully constructed an irrep as 416 presented in § 5. Furthermore, as demonstrated in Lemmas 1, 2 and 3, the joint-equivariance enables 417 a natural handling of (not only fully-connected but any) deep structures. By seemingly making a 418 small adjustment to the conditions, we were able to address a fundamental problem effectively.

419 420

421 422

423

424

425 426

378

379

380 381 382

384

386

387

388

389 390

#### 4 EXAMPLE: DEPTH-*n* JOINT-EQUIVARIANT MACHINE

As pointed out in Lemma 2, the depth-*n* feature map  $\phi_{1:n}$  is joint-*G*-equivariant. Therefore, the following *Y*-valued depth-*n* joint-equivariant machine DLM[ $\gamma; \phi_{1:n}$ ] is  $L^2(X; Y)$ -universal.

**Corollary 1** (Deep Ridgelet Transform). For any maps  $\gamma : X \to \mathbb{C}$  and  $f \in L^2(X; Y)$ , put

$$\mathsf{DLM}[\gamma;\phi_{1:n}](x) := \int_{\Xi_1 \times \dots \times \Xi_n} \gamma(\xi_1,\dots,\xi_n) \phi_n(\bullet,\xi_n) \circ \dots \circ \phi_1(x,\xi_1) \mathrm{d}\boldsymbol{\xi}, \quad x \in X,$$
(18)

427 428

$$\mathbf{R}[f;\psi_{1:n}](\boldsymbol{\xi}) := \int_{\Xi} \langle f(x),\psi_n(\bullet,\xi_n)\circ\cdots\circ\psi_1(x,\xi_n)\rangle_Y \mathrm{d}x, \quad \boldsymbol{\xi}\in\Xi_1\times\cdots\times\Xi_n.$$
(19)

429 430

431 Under the assumptions that  $DLM_{\phi_{1:n}} \circ R_{\psi_{1:n}}$  is bounded, and that  $\pi$  is irreducible, there exists a bilinear form  $((\phi_{1:n}, \psi_{1:n}))$  satisfying  $DLM_{\phi_{1:n}} \circ R_{\psi_{1:n}} = ((\phi_{1:n}, \psi_{1:n})) Id_{L^2(X;Y)}$ .

Again, it extends the original integral representation, and inherits the *linearization* trick of nonlinear parameters  $\boldsymbol{\xi}$  by integrating all the possible parameters (beyond the difference of layers) and indirectly select which parameters to use by weighting on them.

435 436 437

447

453

458

461

466

475 476 477

478 479

481

482 483 484

485

#### 5 EXAMPLE: DEPTH-*n* FULLY-CONNECTED NETWORK

438 439 440 We explain the case of depth-n (precisely, depth-n + 1) fully-connected network. We use the following fact.

**Lemma 4** (Folland (2015, Theorem 7.12)). Let  $\pi_1$  and  $\pi_2$  be representations of locally compact groups  $G_1$  and  $G_2$ , and let  $\pi_1 \otimes \pi_2$  be their outer tensor product, which is a representation of the product group  $G_1 \times G_2$ . Then,  $\pi_1$  and  $\pi_2$  are irreducible if and only if  $\pi_1 \otimes \pi_2$  is irreducible.

444 Set  $X = Y = \mathbb{R}^m$  (input and output domains), and for each  $i \in \{1, ..., n\}$ , set  $X_i := \mathbb{R}^{d_i}$  (with 445  $X_1 = X$  and  $X_{n+1} = Y$ ),  $\Xi_i := \mathbb{R}^{p_i \times d_i} \times \mathbb{R}^{p_i} \times \mathbb{R}^{d_{i+1} \times q_i}$  (parameter domain),  $\sigma_i : \mathbb{R}^{p_i} \to \mathbb{R}^{q_i}$ 446 (activation functions), and define the feature map (vector-valued fully-connected neurons) as

$$\phi_i(\boldsymbol{x}_i, \boldsymbol{\xi}_i) := C_i \sigma_i (A_i \boldsymbol{x}_i - \boldsymbol{b}_i), \quad \boldsymbol{x}_i \in \mathbb{R}^{d_i}, \boldsymbol{\xi}_i = (A_i, \boldsymbol{b}_i, C_i) \in \Xi_i$$
(20)

448 Specifically,  $d_1 = d_{n+1} = m$ . If there is no risk of confusion, we omit writing *i* for simplicity.

450 Let O(m) denote the orthogonal group in dimension m. Let  $G := O(m) \times Aff(m)$  be the product 451 group of O(m) and  $Aff(m) = GL(m) \ltimes \mathbb{R}^m$ . We suppose G acts on the input and output domains 452 as below: For any  $g = (Q, L, t) \in G = O(m) \times (GL(m) \ltimes \mathbb{R}^m)$ ,

$$g \cdot \boldsymbol{x} := L\boldsymbol{x} + \boldsymbol{t}, \ \boldsymbol{x} \in X, \quad \text{and} \quad g \cdot \boldsymbol{y} := v_g[\boldsymbol{y}] := Q\boldsymbol{y}, \ \boldsymbol{y} \in Y.$$
(21)

Namely, the group actions of both O(m) on X and Aff(m) on Y are trivial.

Let  $\pi$  be the induced representation of G on the vector-valued square-integrable functions  $L^2(X;Y)$ , defined by

$$\pi_{g}[\boldsymbol{f}](\boldsymbol{x}) := |\det L|^{-1/2} Q \boldsymbol{f}(L^{-1}(\boldsymbol{x} - \boldsymbol{t})), \quad \boldsymbol{x} \in X, \ \boldsymbol{f} \in L^{2}(X;Y)$$
(22)

459 for each  $g = (Q, L, t) \in O(m) \times (GL(m) \ltimes \mathbb{R}^m)$ .

**Lemma 5.** The above  $\pi : G \to \mathcal{U}(L^2(\mathbb{R}^m; \mathbb{R}^m))$  is irreducible.

462 *Proof.* Recall the representations of O(m) on  $\mathbb{R}^m$  and of Aff(m) on  $L^2(\mathbb{R}^m)$  are respectively ir-463 reducible (see Theorem 3), and  $L^2(\mathbb{R}^m; \mathbb{R}^m)$  is equivalent to the tensor product  $\mathbb{R}^m \otimes L^2(\mathbb{R}^m)$ . 464 Hence by Lemma 4, the representation  $\pi$  of the product group  $O(m) \times Aff(m)$  on the tensor prod-465 uct  $\mathbb{R}^m \otimes L^2(\mathbb{R}^m) = L^2(\mathbb{R}^m; \mathbb{R}^m)$  is irreducible.

Additionally, we put the dual action of G on parameter domain  $\Xi_i$  as below:

$$g \cdot (A_i, \boldsymbol{b}_i, C_i) := \begin{cases} (A_i L^{-1}, \boldsymbol{b}_i + A_i L^{-1} \boldsymbol{t}, C_i), & i = 1\\ (A_i, \boldsymbol{b}_i, C_i), & i \neq 1, n\\ (A_i, \boldsymbol{b}_i, QC_i), & i = n \end{cases}$$
(23)

471 for all  $g = (Q, L, t) \in O(m) \times (GL(m) \ltimes \mathbb{R}^m), \ (A_i, b_i, C_i) \in \Xi_i.$ 

Then, the composition of feature maps  $\phi_{1:n}(x, \xi_{1:n}) := \phi_n(\bullet, \xi_n) \circ \cdots \circ \phi_1(x, \xi_1)$  is joint-*G*-equivariant. In fact,

$$\begin{aligned} \phi_1(g \cdot \boldsymbol{x}, g \cdot \boldsymbol{\xi}_1) &= C_1 \sigma \left( A_1 L^{-1} (L \boldsymbol{x} + \boldsymbol{t}) - (\boldsymbol{b}_1 + A_1 L^{-1} \boldsymbol{t}) \right) = C_1 \sigma (A_1 \boldsymbol{x} - \boldsymbol{b}_1) = \phi_1(\boldsymbol{x}, \boldsymbol{\xi}_1), \\ \phi_i(\boldsymbol{x}, g \cdot \boldsymbol{\xi}_i) &= C_i \sigma (A_i \boldsymbol{x} - \boldsymbol{b}_i) = \phi_i(\boldsymbol{x}, \boldsymbol{\xi}_i), \quad i \neq 1, n \\ \phi_n(\boldsymbol{x}, g \cdot \boldsymbol{\xi}_n) &= Q C_n \sigma (A_n \boldsymbol{x} - \boldsymbol{b}_n) = g \cdot \phi_n(\boldsymbol{x}, \boldsymbol{\xi}_n), \end{aligned}$$

Therefore 
$$\phi_{1:n}(g \cdot \boldsymbol{x}, g \cdot \boldsymbol{\xi}_{1:n}) = g \cdot \phi_{1:n}(\boldsymbol{x}, \boldsymbol{\xi}_{1:n}).$$

480 So by putting depth-n neural network and the corresponding ridgelet transform as below

$$DNN[\gamma; \phi_{1:n}](\boldsymbol{x}) = \int_{\Xi_{1:n}} \gamma(\boldsymbol{\xi}_{1:n}) \phi_{1:n}(\boldsymbol{x}, \boldsymbol{\xi}_{1:n}) d\boldsymbol{\xi}_{1:n},$$
(24)

$$\mathbf{R}[\boldsymbol{f};\psi_{1:n}](\boldsymbol{\xi}_{1:n}) = \int_{\mathbb{R}^m} \boldsymbol{f}(\boldsymbol{x}) \cdot \overline{\psi_{1:n}(\boldsymbol{x},\boldsymbol{\xi}_{1:n})} \mathrm{d}\boldsymbol{x},$$
(25)

Theorem 4 yields the reconstruction formula  $DNN_{\phi_{1:n}} \circ \mathbf{R}_{\psi_{1:n}} = ((\phi_{1:n}, \psi_{1:n})) \operatorname{Id}_{L^2(\mathbb{R}^m;\mathbb{R}^m)}.$ 

## <sup>486</sup> 6 EXAMPLE: QUADRATIC-FORM WITH NONLINEARITY

Here, we present a new network for which the universality was not known.

Let M denote the class of all  $m \times m$ -symmetric matrices equipped with the Lebesgue measure  $dA = \bigwedge_{i>i} da_{ij}$ . Set  $X = \mathbb{R}^m, \Xi = M \times \mathbb{R}^m \times \mathbb{R}$ , and

$$\phi(\boldsymbol{x},\xi) := \sigma(\boldsymbol{x}^{\top}A\boldsymbol{x} + \boldsymbol{x}^{\top}\boldsymbol{b} + c)$$
(26)

for any fixed function  $\sigma : \mathbb{R} \to \mathbb{R}$ . Namely, it is a quadratic-form in x followed by nonlinear activation function  $\sigma$ .

Then, it is joint-invariant with G = Aff(m). In fact, we can put the group actions of  $g = (t, L) \in \mathbb{R}^m \rtimes GL(m)$  on X and  $\Xi$  by

$$(\boldsymbol{t}, \boldsymbol{L}) \cdot \boldsymbol{x} := \boldsymbol{t} + \boldsymbol{L}\boldsymbol{x},\tag{27}$$

$$(\boldsymbol{t}, L) \cdot (A, \boldsymbol{b}, c) := (L^{-\top}AL^{-1}, L^{-\top}\boldsymbol{b} - 2L^{-\top}AL^{-1}\boldsymbol{t}, c + \boldsymbol{t}^{\top}L^{-\top}AL^{-1}\boldsymbol{t} - \boldsymbol{t}^{\top}L^{-\top}\boldsymbol{b}).$$
(28)

Then, the joint-G-action on  $X \times \Xi$  remains the feature map joint-invariant as below.

$$\begin{split} \phi(g \cdot \boldsymbol{x}, g \cdot \boldsymbol{\xi}) &= \sigma((L\boldsymbol{x} + \boldsymbol{t})^\top L^{-\top} A L^{-1} (L\boldsymbol{x} + \boldsymbol{t}) + (L\boldsymbol{x} + \boldsymbol{t})^\top (L^{-\top} \boldsymbol{b} - 2L^{-\top} A L^{-1} \boldsymbol{t}) + ...) \\ &= \sigma(\boldsymbol{x}^\top A \boldsymbol{x} + 2 \boldsymbol{x}^\top A L^{-1} \boldsymbol{t} + \boldsymbol{t}^\top L^{-\top} A L^{-1} \boldsymbol{t} + \\ &+ \boldsymbol{x}^\top \boldsymbol{b} - 2 \boldsymbol{x}^\top A L^{-1} \boldsymbol{t} + \boldsymbol{t}^\top L^{-\top} \boldsymbol{b} - 2 \boldsymbol{t}^\top L^{-\top} A L^{-1} \boldsymbol{t} \\ &+ c + \boldsymbol{t}^\top L^{-\top} A L^{-1} \boldsymbol{t} - \boldsymbol{t}^\top L^{-\top} \boldsymbol{b}) \\ &= \sigma(\boldsymbol{x}^\top A \boldsymbol{x} + \boldsymbol{x}^\top \boldsymbol{b} + c) = \phi(g \cdot \boldsymbol{x}, g \cdot \boldsymbol{\xi}). \end{split}$$

<sup>510</sup> By Theorem 3, the regular representation  $\pi$  of Aff(m) on  $L^2(\mathbb{R}^m)$  is irreducible. Hence as a consequence of the general result, the following network is  $L^2(\mathbb{R}^m)$ -universal.

$$NN[\gamma](\boldsymbol{x}) := \int_{M \times \mathbb{R}^m \times \mathbb{R}} \gamma(A, \boldsymbol{b}, c) \sigma(\boldsymbol{x}^\top A \boldsymbol{x} + \boldsymbol{x}^\top \boldsymbol{b} + c) dA d\boldsymbol{b} dc.$$
(29)

#### 7 DISCUSSION

516 517 518

519

520

521

522

523

524

525

526

527

528 529

530

533

513 514 515

487 488

489

490

491 492 493

496

497

502

> We have developed a systematic method for deriving a ridgelet transform for a wide range of learning machines defined by joint-group-equivariant feature maps, yielding the universal approximation theorems as corollaries. Traditionally, the techniques used in the expressive power analysis of deep networks were different from those used in the analysis of shallow networks, as overviewed in the introduction. Our main theorem unifies the approximation schemes of both deep and shallow networks from the perspective of joint-group-action on the data-parameter domain. Technically, this unification is due to Schur's lemma, a basic and useful result in the representation theory. Thanks to this lemma, the proof of the main theorem is simple, yet the scope of application is wide. The significance of this study lies in revealing the close relationship between machine learning theory and modern algebra. With this study as a catalyst, we expect a major upgrade to machine learning theory from the perspective of modern algebra.

#### References

Andrew R Barron. Universal approximation bounds for superpositions of a sigmoidal function.
 *IEEE Transactions on Information Theory*, 39(3):930–945, 1993.

Michael M. Bronstein, Joan Bruna, Taco Cohen, and Petar Veličković. Geometric Deep Learning:
 Grids, Groups, Graphs, Geodesics, and Gauges. *arXiv preprint: 2104.13478*, 2021.

 536
 537
 538
 Yongqiang Cai. Achieve the Minimum Width of Neural Networks for Universal Approximation. In The Eleventh International Conference on Learning Representations, 2023.

539 Emmanuel Jean Candès. *Ridgelets: theory and applications*. PhD thesis, Standford University, 1998.

540 S. M. Carroll and B. W. Dickinson. Construction of neural nets using the Radon transform. In 541 International Joint Conference on Neural Networks 1989, volume 1, pp. 607–611. IEEE, 1989. 542 Ricky T. O. Chen, Yulia Rubanova, Jesse Bettencourt, and David Duvenaud. Neural Ordinary Dif-543 ferential Equations. In Advances in Neural Information Processing Systems, volume 31, pp. 544 6572–6583, Palais des Congrès de Montréal, Montréal CANADA, 2018. Curran Associates, Inc. 546 Albert Cohen, Ronald DeVore, Guergana Petrova, and Przemysław Wojtaszczyk. Optimal Stable 547 Nonlinear Approximation. Foundations of Computational Mathematics, 22(3):607–648, 2022. 548 Nadav Cohen, Or Sharir, and Amnon Shashua. On the Expressive Power of Deep Learning: A 549 Tensor Analysis. In 29th Annual Conference on Learning Theory, volume 49, pp. 1–31, 2016. 550 551 George Cybenko. Approximation by superpositions of a sigmoidal function. Mathematics of Con-552 trol, Signals, and Systems (MCSS), 2(4):303–314, 1989. 553 I Daubechies, R DeVore, S Foucart, B Hanin, and G Petrova. Nonlinear Approximation and (Deep) 554 ReLU Networks. Constructive Approximation, 55(1):127–172, 2022. 555 556 Weinan E. A Proposal on Machine Learning via Dynamical Systems. Communications in Mathe*matics and Statistics*, 5(1):1–11, 2017. 558 Gerald B. Folland. A Course in Abstract Harmonic Analysis. Chapman and Hall/CRC, New York, 559 second edition, 2015. 560 561 Ken-Ichi Funahashi. On the approximate realization of continuous mappings by neural networks. 562 Neural Networks, 2(3):183-192, 1989. 563 Philipp Grohs, Andreas Klotz, and Felix Voigtlaender. Phase Transitions in Rate Distortion Theory 564 and Deep Learning. Foundations of Computational Mathematics, 23(1):329–392, 2023. 565 566 Eldad Haber and Lars Ruthotto. Stable architectures for deep neural networks. Inverse Problems, 567 34(1):1-22, 2017. 568 Boris Hanin and Mark Sellke. Approximating Continuous Functions by ReLU Nets of Minimal 569 Width. arXiv preprint: 1710.11278, 2017. 570 571 Kurt Hornik, Maxwell Stinchcombe, and Halbert White. Multilayer feedforward networks are uni-572 versal approximators. Neural Networks, 2(5):359–366, 1989. 573 Patrick Kidger and Terry Lyons. Universal Approximation with Deep Narrow Networks. In Pro-574 ceedings of Thirty Third Conference on Learning Theory, volume 125 of Proceedings of Machine 575 Learning Research, pp. 2306–2327. PMLR, 2020. 576 577 Namjun Kim, Chanho Min, and Sejun Park. Minimum width for universal approximation using 578 ReLU networks on compact domain. In The Twelfth International Conference on Learning Rep-579 resentations, 2024. 580 Li'Ang Li, Yifei Duan, Guanghua Ji, and Yongqiang Cai. Minimum Width of Leaky-ReLU Neural 581 Networks for Uniform Universal Approximation. In Proceedings of the 40th International Con-582 ference on Machine Learning, volume 202 of Proceedings of Machine Learning Research, pp. 583 19460-19470. PMLR, 2023. 584 585 Qianxiao Li and Shuji Hao. An Optimal Control Approach to Deep Learning and Applications to Discrete-Weight Neural Networks. In Proceedings of The 35th International Conference on 586 Machine Learning, volume 80, pp. 2985–2994, Stockholm, 2018. PMLR. 588 Hongzhou Lin and Stefanie Jegelka. ResNet with one-neuron hidden layers is a Universal Approxi-589 mator. In Advances in Neural Information Processing Systems, volume 31, Montreal, BC, 2018. 590 Curran Associates, Inc. Zhou Lu, Hongming Pu, Feicheng Wang, Zhiqiang Hu, and Liwei Wang. The Expressive Power 592 of Neural Networks: A View from the Width. In Advances in Neural Information Processing 593

Systems, volume 30. Curran Associates, Inc., 2017.

594 595	Noboru Murata. An integral representation of functions using three-layered networks and their approximation bounds. <i>Neural Networks</i> , 9(6):947–956, 1996.
596 597 598	Sejun Park, Chulhee Yun, Jaeho Lee, and Jinwoo Shin. Minimum Width for Universal Approxima- tion. In International Conference on Learning Representations, 2021.
599 600	Guergana Petrova and Przemyslaw Wojtaszczyk. Limitations on approximation by deep and shallow neural networks. <i>Journal of Machine Learning Research</i> , 24(353):1–38, 2023.
602 603	Jonathan W Siegel. Optimal Approximation Rates for Deep ReLU Neural Networks on Sobolev and Besov Spaces. <i>Journal of Machine Learning Research</i> , 24(357):1–52, 2023.
604 605	Sho Sonoda and Noboru Murata. Neural network with unbounded activation functions is universal approximator. <i>Applied and Computational Harmonic Analysis</i> , 43(2):233–268, 2017a.
607 608 609	Sho Sonoda and Noboru Murata. Transportation analysis of denoising autoencoders: a novel method for analyzing deep neural networks. In <i>NIPS 2017 Workshop on Optimal Transport &amp; Machine Learning (OTML)</i> , pp. 1–10, Long Beach, 2017b.
610 611 612 613	Sho Sonoda, Isao Ishikawa, and Masahiro Ikeda. Ridge Regression with Over-Parametrized Two- Layer Networks Converge to Ridgelet Spectrum. In Proceedings of The 24th International Con- ference on Artificial Intelligence and Statistics (AISTATS) 2021, volume 130, pp. 2674–2682. PMLR, 2021a.
614 615 616	Sho Sonoda, Isao Ishikawa, and Masahiro Ikeda. Ghosts in Neural Networks: Existence, Structure and Role of Infinite-Dimensional Null Space. arXiv preprint: 2106.04770, 2021b.
617 618 619	Sho Sonoda, Isao Ishikawa, and Masahiro Ikeda. Universality of Group Convolutional Neural Networks Based on Ridgelet Analysis on Groups. In Advances in Neural Information Processing Systems 35, pp. 38680–38694, New Orleans, Louisiana, USA, 2022a. Curran Associates, Inc.
620 621 622 623	Sho Sonoda, Isao Ishikawa, and Masahiro Ikeda. Fully-Connected Network on Noncompact Symmetric Space and Ridgelet Transform based on Helgason-Fourier Analysis. In <i>Proceedings of the 39th International Conference on Machine Learning</i> , volume 162, pp. 20405–20422, Baltimore, Maryland, USA, 2022b. PMLR.
625 626 627 628	Sho Sonoda, Hideyuki Ishi, Isao Ishikawa, and Masahiro Ikeda. Joint Group Invariant Functions on Data-Parameter Domain Induce Universal Neural Networks. In Proceedings of the 2nd NeurIPS Workshop on Symmetry and Geometry in Neural Representations, Proceedings of Machine Learning Research, pp. 129–144. PMLR, 2024a.
629 630 631	Sho Sonoda, Isao Ishikawa, and Masahiro Ikeda. A unified Fourier slice method to derive ridgelet transform for a variety of depth-2 neural networks. <i>Journal of Statistical Planning and Inference</i> , 233:106184, 2024b.
632 633 634	Matus Telgarsky. Benefits of depth in neural networks. In 29th Annual Conference on Learning Theory, pp. 1–23, 2016.
635 636 637 638	Hayata Yamasaki, Sathyawageeswar Subramanian, Satoshi Hayakawa, and Sho Sonoda. Quantum Ridgelet Transform: Winning Lottery Ticket of Neural Networks with Quantum Computation. In Proceedings of the 40th International Conference on Machine Learning, volume 202 of Proceed- ings of Machine Learning Research, pp. 39008–39034, Honolulu, Hawaii, USA, 2023. PMLR.
639 640 641	Dmitry Yarotsky. Error bounds for approximations with deep ReLU networks. <i>Neural Networks</i> , 94:103–114, 2017.
642 643 644	Dmitry Yarotsky. Optimal approximation of continuous functions by very deep ReLU networks. In <i>Proceedings of the 31st Conference On Learning Theory</i> , volume 75 of <i>Proceedings of Machine Learning Research</i> , pp. 639–649. PMLR, 2018.
645 646 647	Dmitry Yarotsky and Anton Zhevnerchuk. The phase diagram of approximation rates for deep neural networks. In <i>Advances in Neural Information Processing Systems</i> , volume 33, pp. 13005–13015. Curran Associates, Inc., 2020.

А

 $L^2(X;Y),$ 

$$\pi_g[\pi_h[f]](x) = g \cdot (h \cdot f(h^{-1} \cdot (g^{-1} \cdot x))) = (gh) \cdot f((gh)^{-1} \cdot x) = \pi_{gh}[f](x),$$

*Proof.* Recall that the representation v of G on Y is unitary. So, for any  $g, h \in G$  and  $f \in G$ 

and for any  $g \in G$  and  $f_1, f_2 \in L^2(X; Y)$ ,

A.1 UNITARITY OF REPRESENTATIONS

**Lemma 6.**  $\pi$  is a unitary representation of G on  $L^2(X;Y)$ .

PROOFS

$$\langle \pi_g[f_1], \pi_g[f_2] \rangle_{L^2(X;Y)} = \int_X \langle \upsilon_g[f_1(g^{-1} \cdot x)], \upsilon_g[f_2(g^{-1} \cdot x)] \rangle_Y \mathrm{d}x$$
  
= 
$$\int_X \langle f_1(x), \upsilon_g^*[\upsilon_g[f_2(x)]] \rangle_Y \mathrm{d}x = \langle f_1, f_2 \rangle_{L^2(X;Y)}. \qquad \Box$$

**Lemma 7.**  $\hat{\pi}$  is a unitary representation of G on  $L^2(\Xi)$ .

*Proof.* For any  $g, h \in G$  and  $\gamma \in L^2(\Xi)$ ,

$$\widehat{\pi}_g[\widehat{\pi}_h[\gamma]](\xi) = \gamma(h^{-1} \cdot (g^{-1} \cdot \xi)) = \gamma((gh)^{-1} \cdot \xi) = \widehat{\pi}_{gh}[f](x),$$

and for any  $g \in G$  and  $\gamma_1, \gamma_2 \in L^2(\Xi)$ ,

$$\langle \hat{\pi}_{g}[\gamma_{1}], \hat{\pi}_{g}[\gamma_{2}] \rangle_{L^{2}(\Xi)} = \int_{\Xi} \gamma_{1}(g^{-1} \cdot \xi) \overline{\gamma_{2}(g^{-1} \cdot \xi)} \mathrm{d}\xi$$
$$= \int_{\Xi} \gamma_{1}(\xi) \overline{\gamma_{2}(\xi)} \mathrm{d}\xi = \langle \gamma_{1}, \gamma_{2} \rangle_{L^{2}(\Xi)}.$$

#### A.2 PROOF OF LEMMA 2

Proof. For any  $g \in G, x \in X$ , and  $\xi_{1:n} \in \Xi_{1:n}$ , we have  $\begin{aligned} \phi_{1:n}(g \cdot x, g \cdot \xi_{1:n}) &= \phi_n(\bullet, g \cdot \xi_n) \circ \cdots \circ \phi_2(\bullet, g \cdot \xi_2) \circ \phi_1(g \cdot x, g \cdot \xi_1) \\ &= \phi_n(\bullet, g \cdot \xi_n) \circ \cdots \circ \phi_2(g \cdot \bullet, g \cdot \xi_2) \circ \phi_1(x, \xi_1) \\ &\vdots \\ &= \phi_n(g \cdot \bullet, g \cdot \xi_n) \circ \cdots \circ \phi_2(\bullet, \xi_2) \circ \phi_1(x, \xi_1) \\ &= g \cdot \phi_n(\bullet, \xi_n) \circ \cdots \circ \phi_2(\bullet, \xi_2) \circ \phi_1(x, \xi_1) \\ &= g \cdot \phi_{1:n}(x, \xi_{1:n}). \end{aligned}$ 

# 702 B EXAMPLE: DEPTH-*n* GROUP CONVOLUTIONAL NETWORK

As mentioned in Remark 1, all the classical equivariant feature maps are automatically jointequivariant. So, once the irreducibility of representation  $\pi$  is verified, our main theorem can state the ridgelet transform for classical equivariant networks. Here, we explain another usage of this study to present the ridgelet transform for *depth-n* group convolutional networks (GCNs).

<sup>708</sup> In a previous study (Sonoda et al., 2022a), the authors have introduced a general formulation of <sup>709</sup> *depth-2* GCNs, covering a wide range of typical group equivariant networks such as DeepSets and <sup>710</sup> E(n)-equivariant maps, and presented the ridgelet transform for them in a unified manner. The <sup>711</sup> ridgelet transform for GCNs was derived and shown to satisfy the reconstruction formula based on <sup>712</sup> the *Fourier expression method* (see Sonoda et al., 2024b, for more details), another proof technique <sup>713</sup> for ridgelet transforms that does not require the irreducibility assumption but is limited to depth-2 <sup>714</sup> learning machines.

<sup>715</sup> In the following, we extend the ridgelet transform for GCNs from depth-2 to *depth-n* by reviewing it <sup>716</sup> from the group theoretic perspective. The main idea is to turn the depth-*n* fully-connected network <sup>717</sup> (FCN)  $\phi_{1:n}$  in § 5 to a depth-*n G*-convolutional network  $\psi_{1:n}$  by following the construction of the <sup>718</sup> previous study.

719

740

747 748 749

750

751 752

753

755

720 B.1 NOTATIONS 721

To begin with, we introduce the affine group  $A := Aff(m) = GL(m) \ltimes \mathbb{R}^m$  as an auxiliary group, besides the primary group G in consideration. Eventually, the irreducibility assumption is required not for G but for A. We note that A and G need not satisfy the inclusion relations neither  $A \leq G$ nor  $G \leq A$ . In accordance with the previous study, we write  $T_g[\bullet]$  for G-action, and write  $\alpha \cdot \bullet$  for A-action if needed. Caution: Different from § 5, we assign A (instead of G) for the affine group acting on  $\phi_{1:n}$ . So, we will turn a joint-A-equivariant map  $\phi_{1:n}$  to G-equivariant map  $\psi_{1:n}$ .

Suppose G acts on each intermediate feature space  $X_i = \mathbb{R}^{d_i}$  (and acts trivially on each parameter space  $\Xi_i$ ). By abusing notations, we use the same symbols T for every G-actions on  $X_i$ , and  $\tau$  for every G-actions on function  $f : X \to Y$ , where X and Y denote any G-sets such as G itself and  $X_i$ , defined by  $\tau_g[f](x) := T_g[f(T_{g^{-1}}[x])]$  for every  $g, h \in G$ . By  $L^2_G(X;Y)$ , we denote the space of G-equivariant Y-valued functions f on X that is square-integrable at the identity element  $1_G$  of G, namely  $L^2_G(X;Y) = \{f \in \hom_G(X,Y^G) \mid ||f(\bullet)(1_G)||_{L^2(X;Y)} < \infty\} \cong \{\tau_{\bullet}[f_1] \mid f_1 \in L^2(X;Y)\}$ .

#### 735 736 B.2 *G*-Convolutional Feature Map

From the *i*-th layer fully-connected map  $\phi_i : X_i \times \Xi_i \to X_{i+1}$ , we define the *i*-th layer *G*convolutional map  $\psi_i : X_i \times \Xi_i \to X_{i+1}^G$  as follows: For all  $\boldsymbol{x}_i \in X_i$  and  $\boldsymbol{\xi}_i = (A_i, \boldsymbol{b}_i, C_i) \in \Xi_i$ ,

$$\psi_i(\boldsymbol{x}_i, \boldsymbol{\xi}_i)(g) := \tau_g[\phi_i](\boldsymbol{x}_i, \boldsymbol{\xi}_i) = T_g[(C_i \sigma_i (A_i T_{g^{-1}}[\boldsymbol{x}_i] - \boldsymbol{b}_i)], \quad g \in G.$$
(30)

This is called a *G*-convolutional map because by appropriately specifying the *G*-action *T*, the expression  $A_i T_{g^{-1}}[x_i]$  can reproduce a variety of *G*-convolution product, say  $a *_T x$ , employed in such as DeepSets and E(n)-equivariant maps (see Section 5 of Sonoda et al., 2022a).

Similarly to Lemma 1, each G-convolutional map  $\psi_i$  is G-equivariant in the classical sense (or joint-G-equivariant with trivial G-action on parameters  $\xi_i$ ) because for any  $g, h \in G$ ,

$$\psi_i(T_g[\boldsymbol{x}_i], \boldsymbol{\xi}_i)(h) = T_h[\phi_i(T_{h^{-1}}[T_g[\boldsymbol{x}_i]], \boldsymbol{\xi}_i)] = T_g[T_{g^{-1}h}[\phi_i(T_{(g^{-1}h)^{-1}}[\boldsymbol{x}_i], \boldsymbol{\xi}_i)]] = \tau_g[\psi_i(\boldsymbol{x}_i, \boldsymbol{\xi}_i)](h).$$
(31)

We note that the G-equivariance holds for any activation function  $\sigma_i$ , because it is applied in the element-wise manner in G.

#### B.3 G-CONVOLUTIONAL NETWORK AND RIDGELET TRANSFORM

Next, we define the depth-*n G*-convolutional map  $\psi_{1:n}: X \times \Xi_{1:n} \to Y^G$  by their compositions:

$$\psi_{1:n}(\boldsymbol{x},\boldsymbol{\xi}_{1:n})(g) := \psi_n(\bullet,\boldsymbol{\xi}_n)(g) \circ \cdots \circ \psi_1(\boldsymbol{x},\boldsymbol{\xi}_1)(g), \tag{32}$$

and define the depth-n G-convolutional network by its integration:

$$\operatorname{GCN}[\gamma](\boldsymbol{x})(g) := \int_{\Xi_{1:n}} \gamma(\boldsymbol{\xi}_{1:n}) \psi_{1:n}(\boldsymbol{x}, \boldsymbol{\xi}_{1:n})(g) \mathrm{d}\boldsymbol{\xi}_{1:n}.$$
(33)

Recall that the depth-n FCN and its ridgelet transform,  $R_{fc}$ , are given in (24) and (25) as below.

$$DNN[\gamma](\boldsymbol{x}) := \int_{\Xi_{1:n}} \gamma(\boldsymbol{\xi}_{1:n}) \phi_{1:n}(\boldsymbol{x}, \boldsymbol{\xi}_{1:n}) \mathrm{d}\boldsymbol{\xi}_{1:n},$$
(34)

$$R_{fc}[f](\boldsymbol{\xi}_{1:n}) = \int_{\mathbb{R}^m} \boldsymbol{f}(\boldsymbol{x}) \cdot \overline{\phi'_{1:n}(\boldsymbol{x}, \boldsymbol{\xi}_{1:n})} d\boldsymbol{x}.$$
(35)

Then, as a consequence of Lemmas 2 and 3, we have the following.

768 Lemma 8.  $GCN[\gamma](\boldsymbol{x})(g) = \tau_g[DNN[\gamma]](\boldsymbol{x}).$ 

Proof. In fact,

$$\psi_{1:n}(\boldsymbol{x}, \boldsymbol{\xi}_{1:n})(g) = T_g[\phi_n(\bullet, \boldsymbol{\xi}_n) \circ \cdots \circ \phi_1(T_{g^{-1}}[\boldsymbol{x}], \boldsymbol{\xi}_1)] \\ = T_g[\phi_{1:n}(T_{g^{-1}}[\boldsymbol{x}], \boldsymbol{\xi}_{1:n})] = \tau_g[\phi_{1:n}](\boldsymbol{x}, \boldsymbol{\xi}_{1:n}),$$

and thus

$$\operatorname{GCN}[\gamma](\boldsymbol{x})(g) = \int_{\Xi_{1:n}} \gamma(\boldsymbol{\xi}_{1:n}) \tau_g[\phi_{1:n}](\boldsymbol{x}, \boldsymbol{\xi}_{1:n}) \mathrm{d}\boldsymbol{\xi}_{1:n} = \tau_g[\operatorname{DNN}[\gamma]](\boldsymbol{x}).$$

Finally, the ridgelet transform  $R_{conv}$  for depth-*n* GCNs is given by using  $R_{fc}$  for depth-*n* FCNs as below: For any  $f \in L^2_G(X : Y)$ , put

$$\mathbf{R}_{\operatorname{conv}}[f](\boldsymbol{\xi}_{1:n}) := \mathbf{R}_{\operatorname{fc}}[f(\bullet)(1_G)](\boldsymbol{\xi}_{1:n}).$$
(36)

782  
783 Proof. 
$$GCN[R_{conv}[f]](x)(g) = \tau_g[DNN[R_{fc}[f(\bullet)(1_G)]]](x) = \tau_g[f(\bullet)(1_G)](x) = f(x)(g).$$

<sup>784</sup> In other words, the ridgelet transform only encodes the information of function f at  $1_G$  because the <sup>785</sup> essential information of f is summarized at  $1_G$  due to its G-equivariance, and the G-convolutions <sup>786</sup> in depth-n GCN have a mechanism to automatically expand the summarized information to entire <sup>787</sup> G by using G-equivariance. When the depth n = 2, it reduces to the ridgelet transform for depth-2 <sup>788</sup> GCNs presented in Theorem 1 of Sonoda et al. (2022a).

We remark that the base feature map  $\phi_i$  need not be the fully-connected network. As suggested from the construction at (30), it can be an arbitrary joint-*A*-equivariant map. However, we need to verify the irreducibility of representation  $\pi$  of *A* on  $L^2(X;Y)$  in general.