
Constructive Universal Approximation Theorems for Deep Joint-Equivariant Networks by Schur’s Lemma

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 We present a unified constructive universal approximation theorem covering a wide
2 range of learning machines including both shallow and deep neural networks based
3 on the group representation theory. Constructive here means that the distribution
4 of parameters is given in a closed-form expression (called the *ridgelet transform*).
5 Contrary to the case of shallow models, expressive power analysis of deep models
6 has been conducted in a case-by-case manner. Recently, Sonoda et al. [33, 32]
7 developed a systematic method to show a constructive approximation theorem
8 from *scalar-valued joint-group-invariant* feature maps, covering a formal deep
9 network. However, each hidden layer was formalized as an abstract group action, so
10 it was not possible to cover real deep networks defined by composites of nonlinear
11 activation function. In this study, we extend the method for *vector-valued joint-*
12 *group-equivariant* feature maps, so to cover such real networks.

13 1 Introduction

14 An ultimate goal of the deep learning theory is to characterize the internal data processing procedure
15 inside deep neural networks obtained by deep learning. We may formulate this problem as a functional
16 equation problem: Let \mathcal{F} denote a class of data generating functions, and let $\text{DNN}[\gamma]$ denote a certain
17 deep neural network with parameter γ . Given a function $f \in \mathcal{F}$, find an unknown parameter γ so that
18 network $\text{DNN}[\gamma]$ represents function f , i.e.

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

19 which we call a *DNN equation*. An ordinary learning problem by empirical risk minimization, such
20 as minimizing $\sum_{i=1}^n |\text{DNN}[\gamma](x_i) - f(x_i)|^2$ with respect to γ , is understood as a weak form (or a
21 variational form) of this equation. Therefore, characterizing the solution space of this equation leads
22 to understanding the parameters obtained by deep learning. Following previous studies [21, 3, 28–
23 31], we call a solution operator R that satisfies $\text{DNN}[R[f]] = f$ a *ridgelet transform*. Once such a
24 solution operator R is found, we can conclude a *universality* of the DNN in consideration because the
25 reconstruction formula $\text{DNN}[R[f]] = f$ implies for any $f \in \mathcal{F}$ there exists a DNN that represents f .
26 In particular, when $R[f]$ is found in a closed-form manner, then it leads to a *constructive* proof of the
27 universality since $R[f]$ could indicate how to assign parameters.

28 When the network has only one infinitely-wide hidden layer, though it is not deep but shallow, the
29 characterization problem has been well investigated. For example, the learning dynamics and the
30 global convergence property (of SGD) are well studied in the mean field theory [22, 25, 20, 5] and the
31 Langevin dynamics theory [35], and even closed-form solution operator to a “shallow” NN equation,
32 the original ridgelet transform, has already been presented [28–31].

33 On the other hand, when the network has more than one hidden layer, the problem is far from
34 solved, and it is common to either consider infinitely-deep mathematical models such as Neural

35 ODEs [27, 9, 17, 12, 4], or handcraft inner feature maps depending on the problem. For example,
 36 construction methods such as the Telgarsky sawtooth function (or the Yarotsky scheme) and bit
 37 extraction techniques [7, 36–39, 8, 6, 26, 24, 11] have been developed to demonstrate the depth
 38 separation, super-convergence, and minmax optimality of deep ReLU networks. Various feature maps
 39 have also been handcrafted in the contexts of geometric deep learning [1] and deep narrow networks
 40 [19, 13, 18, 14, 23, 16, 2, 15]. Needless to say, there is no guarantee that these handcrafted feature
 41 maps are acquired by deep learning, so these analyses are considered to be analyses of possible
 42 worlds.

43 Recently, Sonoda et al. [33, 32] discovered a rich class of ridgelet transforms for learning machines
 44 defined by *scalar-valued joint-group-invariant* feature maps, covering both depth-2 fully-connected
 45 networks and the formal deep network (FDN), yielding the first ridgelet transform for deep models.
 46 Their theory is indeed a breakthrough because it could cover both deep and shallow models simulta-
 47 neously. However, each hidden layer in the FDN has to be formalized as an abstract group action,
 48 so it was not possible to cover real deep networks defined by composites of nonlinear activation
 49 function. In this study, we extend their arguments for *vector-valued joint-group-equivariant* feature
 50 maps (Theorem 3 and Corollary 1), so to cover such real networks. As an important example, in
 51 § 4.2, we obtained the ridgelet transform for a more realistic DNN, the depth- n fully-connected
 52 network with an arbitrary activation function (not limited to ReLU), without handcrafting network
 53 architecture. In other words, it is a constructive proof of the $L^2(\mathbb{R}^m; \mathbb{R}^m)$ -universality of the DNNs,
 54 and an explicit characterization of the solution space of the DNN equation for more realistic setup.

55 Thanks to Schur’s lemma, a basic and useful result in the representation theory, the proof of the main
 56 theorem is surprisingly simple, yet the scope of application is wide. The significance of this study
 57 lies in revealing the close relationship between machine learning theory and modern algebra. With
 58 this study as a catalyst, we expect a major upgrade to machine learning theory from the perspective
 59 of modern algebra.

60 2 Preliminaries

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

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

69 2.1 Integral Representation and Ridgelet Transform for Depth-2 Fully-Connected Network

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

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

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

72 The integration over all the hidden parameters $(\mathbf{a}, b) \in \mathbb{R}^m \times \mathbb{R}$ means all the neurons $\{\mathbf{x} \mapsto$
 73 $\sigma(\mathbf{a} \cdot \mathbf{x} - b) \mid (\mathbf{a}, b) \in \mathbb{R}^m \times \mathbb{R}\}$ are summed (or integrated, to be precise) with weight γ , hence
 74 formally $S_\sigma[\gamma]$ is understood as a continuous neural network with a single hidden layer. We note,
 75 however, when γ is a finite sum of point measures such as $\gamma_p = \sum_{i=1}^p c_i \delta_{(\mathbf{a}_i, b_i)}$ (by appropriately
 76 extending the class of γ to Borel measures), then it can also reproduce a finite width network

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

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

79 Next, we introduce the ridgelet transform, which is known to be a right-inverse operator to S_σ .

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

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

81 We call R_ρ a ridgelet transform.

82 To be precise, it satisfies the following reconstruction formula.

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

$$S_\sigma \circ R_\rho[f] = ((\sigma, \rho))f, \quad (5)$$

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

87 See Sonoda et al. [29, Theorem 6] for the proof. In particular, according to Sonoda et al. [29,
 88 Lemma 9], for any activation function σ , there always exists ρ satisfying $((\sigma, \rho)) = 1$. Here, σ
 89 being a tempered distribution means that typical activation functions are covered such as ReLU, step
 90 function, tanh, gaussian, etc... We can interpret the reconstruction formula as a universality theorem
 91 of continuous neural networks, since for any given data generating function f , a network with output
 92 weight $\gamma_f = R_\rho[f]$ reproduces f (up to factor $((\sigma, \rho))$), i.e. $S[\gamma_f] = f$. In other words, the ridgelet
 93 transform indicates how the network parameters should be organized so that the network represents
 94 an individual function f .

95 The original ridgelet transform was discovered by Murata [21] and Candès [3]. It is recently extended
 96 to a few modern networks by the Fourier slice method [34, see e.g.]. In this study, we present a
 97 systematic scheme to find the ridgelet transform for a variety of given network architecture based on
 98 the group theoretic arguments.

99 2.2 Irreducible Unitary Representation and Schur's Lemma

100 Let G be a locally compact group, \mathcal{H} be a nonzero Hilbert space, and $\mathcal{U}(\mathcal{H})$ be the group of unitary
 101 operators on \mathcal{H} . For example, any finite group, discrete group, compact group, and finite-dimensional
 102 Lie group are locally compact, while an infinite-dimensional Lie group is not locally compact. A
 103 *unitary representation* π of G on \mathcal{H} is a group homomorphism that is continuous with respect to
 104 the strong operator topology—that is, a map $\pi : G \rightarrow \mathcal{U}(\mathcal{H})$ satisfying $\pi(gh) = \pi(g)\pi(h)$ and
 105 $\pi(g^{-1}) = \pi(g)^{-1}$, and for any $\psi \in \mathcal{H}$, the map $G \ni g \mapsto \pi(g)[\psi] \in \mathcal{H}$ is continuous.

106 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
 107 $g \in G$. Particularly, π is called *irreducible* when it does not admit any nontrivial invariant subspace
 108 $\mathcal{M} \neq \{0\}$ nor \mathcal{H} . The following theorem is a fundamental result of group representation theory that
 109 characterizes the irreducibility.

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

113 See Folland [10, Theorem 3.5(a)] for the proof. We use this as a key step in the proof of our main
 114 theorem.

115 2.3 Calculus on Locally Compact Group

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

$$\int_G f(sg) d_l g = \int_G f(g) d_l g, \quad \text{and} \quad \int_G f(gs) d_r g = \int_G f(g) d_r g.$$

118 Let X be a G -space with transitive left (resp. right) G -action $g \cdot x$ (resp. $x \cdot g$) for any $(g, x) \in G \times X$.
 119 Then, we can further induce the left (resp. right) invariant measure $d_l x$ (resp. $d_r x$) so that for any
 120 $f \in C_c(G)$,

$$\int_X f(x) d_l x := \int_G f(g \cdot o) d_l g, \quad \text{resp.} \quad \int_X f(x) d_r x := \int_G f(o \cdot g) d_r g,$$

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

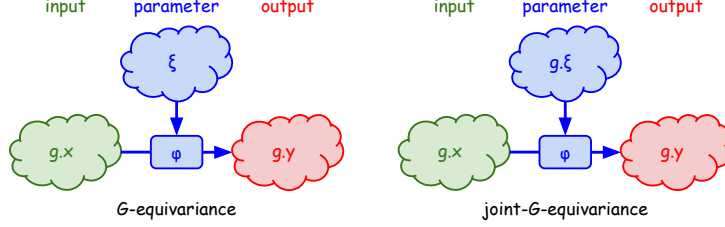


Figure 1: An ordinary G -equivariant feature map $\phi : X \times \Xi \rightarrow Y$ is a subclass of joint- G -equivariant map where the G -action on parameter domain Ξ is *trivial*, i.e. $g \cdot \xi = \xi$

122 3 Main Results

123 We introduce the joint-group-equivariant feature map, and present the ridgelet transforms for learning
 124 machines defined by joint-group-equivariant feature maps, yielding the universality of deep models.

125 Let G be a locally compact group equipped with a left invariant measure dg . Let X and Ξ be
 126 G -spaces equipped with G -invariant measures dx and $d\xi$, called the data domain and the parameter
 127 domain, respectively. Particularly, we call the product space $X \times \Xi$ the *data-parameter domain* (like
 128 time-frequency domain), and call any map ϕ on data-parameter domain $X \times \Xi$ a *feature map*. Let \mathcal{H}
 129 be a separable Hilbert space, let $\mathcal{U}(\mathcal{H})$ be the space of unitary operators on \mathcal{H} , and let $v : G \rightarrow \mathcal{U}(\mathcal{H})$
 130 be a unitary representation of G on \mathcal{H} . If there is no danger of confusion, we use the same symbol \cdot
 131 for the G -actions on X , \mathcal{H} , and Ξ (e.g., $g \cdot x$, $g \cdot v$, and $g \cdot \xi$).

132 In the main theorem, the irreducibility of the following unitary representation π will be a sufficient
 133 condition for the universality. Let $L^2(X; \mathcal{H})$ denote the space of \mathcal{H} -valued square-integrable functions
 134 on X equipped with the inner product $\langle \phi, \psi \rangle_{L^2(X; \mathcal{H})} := \int_X \langle \phi(x), \psi(x) \rangle_{\mathcal{H}} dx$. Put

$$\pi_g[f](x) := g \cdot f(g^{-1} \cdot x), \quad x \in X, f \in L^2(X; \mathcal{H}), g \in G. \quad (6)$$

135 Then, it is a unitary representation of G on $L^2(X; \mathcal{H})$. In fact, $\pi_g[\pi_h[f]](x) = g \cdot h \cdot f(h^{-1} \cdot g^{-1} \cdot x) =$
 136 $(gh) \cdot f((gh)^{-1} \cdot x) = \pi_{gh}[f](x)$, and $\langle \pi_g[f_1], \pi_g[f_2] \rangle_{L^2(X; \mathcal{H})} = \int_X \langle v_g[f_1](g^{-1} \cdot x), v_g[f_2](g^{-1} \cdot$
 137 $x) \rangle_{\mathcal{H}} dx = \int_X \langle f_1(x), v_g^*[v_g[f_2]](x) \rangle_{\mathcal{H}} dx = \langle f_1, f_2 \rangle_{L^2(X; \mathcal{H})}$.

138 In addition, let $L^2(\Xi)$ denote the space of \mathbb{C} -valued square-integrable functions on Ξ , and let $\hat{\pi}$ be
 139 the left-regular representation of G on $L^2(\Xi)$ given by

$$\hat{\pi}_g[\gamma](\xi) := \gamma(g^{-1} \cdot \xi), \quad \xi \in \Xi, \gamma \in L^2(\Xi), g \in G. \quad (7)$$

140 Similarly to π , $\hat{\pi}$ is also a unitary representation.

141 **Definition 3 (Joint G -Equivariant Feature Map).** Let X, Y be data domains, and Ξ be a parameter
 142 domain (with G -actions). We say a feature map $\phi : X \times \Xi \rightarrow Y$ is *joint- G -equivariant* when

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

143 holds for all $g \in G$. In other words, ϕ is a homomorphism (or G -map) of G -sets from $X \times \Xi$ to
 144 Y . So by $\text{hom}_G(X \times \Xi, Y)$, we denote the collection of all joint- G -equivariant maps. Additionally,
 145 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*.

146 *Remark 1.* The joint- G -equivariance extends an ordinary notion of G -equivariance, i.e. $\phi(g \cdot x, \xi) =$
 147 $g \cdot \phi(x, \xi)$. In fact, G -equivariance is a special case of joint- G -equivariance where G acts trivially on
 148 parameter domain, i.e. $g \cdot \xi = \xi$ (see also Figure 1).

149 In order to construct a (non-joint) group-equivariant network, we must carefully and precisely design
 150 the network architecture [see, e.g., a textbook of geometric deep learning 1]. On the other hand, we
 151 can easily and systematically construct joint- G -equivariant network from (not at all equivariant but)
 152 any map $f : X \rightarrow Y$ according to the following Lemmas 1 and 2.

153 **Lemma 1.** Suppose group G acts on sets X and Y . Fix an arbitrary map $f : X \rightarrow Y$, and put
 154 $\phi(x, g) := g \cdot f(g^{-1} \cdot x)$ for every $x \in X$ and $g \in G$. Then, $\phi : X \times G \rightarrow Y$ is joint- G -equivariant.

155 *Proof.* Straightforward. For any $g \in G$, $\phi(g \cdot x, g \cdot h) = (gh) \cdot f((gh)^{-1} \cdot (g \cdot x)) = g \cdot \phi(x, h)$. \square

156 **Lemma 2** (Depth- n Feature Map $\phi_{1:n}$). *Given a sequence of G -equivariant feature maps $\phi_i : X_{i-1} \times \Xi_i \rightarrow X_i$ ($i = 1, \dots, n$), put $\phi_{1:n} : X_0 \times \Xi_1 \times \dots \times \Xi_n \rightarrow X_n$ by*

$$\phi_{1:n}(x, \xi_1, \dots, \xi_n) := \phi_n(\bullet, \xi_n) \circ \dots \circ \phi_1(x, \xi_1). \quad (9)$$

158 *Then, $\phi_{1:n}$ is G -equivariant. Following the custom of counting the number of parameter domains $(\Xi_i)_{i=1}^n$, we say $\phi_{1:n}$ is depth- n .*

160 *Proof.* In fact,

$$\begin{aligned} \phi_{1:n}(g \cdot x, g \cdot \xi_1, \dots, g \cdot \xi_n) &= \phi_n(\bullet, g \cdot \xi_n) \circ \dots \circ \phi_2(\bullet, g \cdot \xi_2) \circ \phi_1(g \cdot x, g \cdot \xi_1) \\ &= \phi_n(\bullet, g \cdot \xi_n) \circ \dots \circ \phi_2(g \cdot \bullet, g \cdot \xi_2) \circ \phi_1(x, \xi_1) \\ &\quad \vdots \\ &= \phi_n(g \cdot \bullet, g \cdot \xi_n) \circ \dots \circ \phi_2(\bullet, \xi_2) \circ \phi_1(x, \xi_1) \\ &= g \cdot \phi_n(\bullet, \xi_n) \circ \dots \circ \phi_2(\bullet, \xi_2) \circ \phi_1(x, \xi_1) \\ &= g \cdot \phi_{1:n}(x, \xi_1, \dots, \xi_n). \quad \square \end{aligned}$$

161 **Definition 4** (ϕ -Network). For any vector-valued map $\phi : X \times \Xi \rightarrow \mathcal{H}$ and scalar-valued map $\gamma : \Xi \rightarrow \mathbb{C}$, define a vector-valued map $X \rightarrow \mathcal{H}$ by

$$\text{NN}[\gamma; \phi](x) := \int_{\Xi} \gamma(\xi) \phi(x, \xi) d\xi, \quad x \in X, \quad (10)$$

163 where the integral is understood as the Bocher integral.

164 We call the integral transform $\text{NN}[\bullet; \phi]$ a ϕ -transform, and each individual image $\text{NN}[\gamma; \phi]$ a ϕ -network for short. The ϕ -network extends the original integral representation. In particular, it inherits the concept of integrating all the possible parameters ξ and indirectly select which parameters to use by weighting on them, which *linearize* parametrization by lifting nonlinear parameters ξ to linear parameter γ .

169 **Definition 5** (ψ -Ridgelet Transform). For any \mathcal{H} -valued feature map $\psi : X \times \Xi \rightarrow \mathcal{H}$ and \mathcal{H} -valued Borel measurable function f on X , put a scalar-valued integral transform

$$\text{R}[f; \psi](\xi) := \int_X \langle f(x), \psi(x, \xi) \rangle_{\mathcal{H}} dx, \quad \xi \in \Xi. \quad (11)$$

171 We call the integral transform $\text{R}[\bullet; \psi]$ a ψ -ridgelet transform for short.

172 As long as the integrals are convergent, ϕ -ridgelet transform is the dual operator of ϕ -transform, since

$$\langle \gamma, \text{R}[f; \psi] \rangle_{L^2(\Xi)} = \int_{X \times \Xi} \gamma(\xi) \langle \phi(x, \xi), f(x) \rangle_{\mathcal{H}} dx d\xi = \langle \text{NN}[\gamma; \phi], f \rangle_{L^2(X; \mathcal{H})}. \quad (12)$$

173 **Theorem 3** (Reconstruction Formula). *Assume (1) \mathcal{H} -valued feature maps $\phi, \psi : X \times \Xi \rightarrow \mathcal{H}$ are joint- G -equivariant, (2) composite operator $\text{NN}_{\phi} \circ \text{R}_{\psi} : L^2(X; \mathcal{H}) \rightarrow L^2(X; \mathcal{H})$ is bounded (i.e., Lipschitz continuous), and (3) the unitary representation π defined in (6) is irreducible. Then, there exists a bilinear form $((\phi, \psi)) \in \mathbb{C}$ (independent of f) such that for any \mathcal{H} -valued square-integrable function $f \in L^2(X; \mathcal{H})$,*

$$\text{NN}_{\phi} \circ \text{R}_{\psi}[f] = ((\phi, \psi)) f. \quad (13)$$

178 In other words, the ψ -ridgelet transform R_{ψ} is a right inverse operator of ϕ -transform NN_{ϕ} as long as $((\phi, \psi)) \neq 0, \infty$.

180 *Proof.* We write $\text{NN}[\bullet; \phi]$ as NN_{ϕ} and $\text{R}[\bullet; \psi]$ as R_{ψ} for short. By using the unitarity of representation $\nu : G \rightarrow \mathcal{U}(\mathcal{H})$, left-invariance of measure dx , and G -equivariance of feature map ψ , for all $g \in G$, we have

$$\begin{aligned} \text{R}_{\psi}[\pi_g[f]](\xi) &= \int_X \langle g \cdot f(g^{-1} \cdot x), \psi(x, \xi) \rangle_{\mathcal{H}} dx = \int_X \langle f(x), g^{-1} \cdot \psi(g \cdot x, \xi) \rangle_{\mathcal{H}} dx \\ &= \int_X \langle f(x), \psi(x, g^{-1} \cdot \xi) \rangle_{\mathcal{H}} dx = \widehat{\pi}_g[\text{R}_{\psi}[f]](\xi). \end{aligned} \quad (14)$$

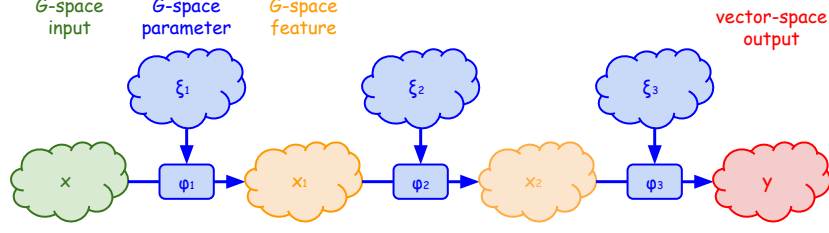


Figure 2: Deep \mathcal{H} -valued joint- G -equivariant network on G -space X is $L^2(X; \mathcal{H})$ -universal when unitary representation π of G on $L^2(X; \mathcal{H})$ is irreducible, and the distribution of parameters for the network to represent a given map $f : X \rightarrow \mathcal{H}$ is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

183 Similarly,

$$\begin{aligned} \text{NN}_\phi[\widehat{\pi}_g[\gamma]](x) &= \int_{\Xi} \gamma(g^{-1} \cdot \xi) \phi(x, \xi) d\xi = \int_{\Xi} \gamma(\xi) \phi(x, g \cdot \xi) d\xi \\ &= \int_{\Xi} \gamma(\xi) (g \cdot \phi(g^{-1} \cdot x, \xi)) d\xi = \pi_g[\text{NN}_\phi[\gamma]](x). \end{aligned} \quad (15)$$

184 Here, $\widehat{\pi}^*$ denotes the dual representation of $\widehat{\pi}$ with respect to $L^2(\Xi)$.

185 As a consequence, $\text{NN}_\phi \circ \mathbb{R}_\psi : L^2(X; \mathcal{H}) \rightarrow L^2(X; \mathcal{H})$ commutes with π as below

$$\text{NN}_\phi \circ \mathbb{R}_\psi \circ \pi_g = \text{NN}_\phi \circ \widehat{\pi}_g \circ \mathbb{R}_\psi = \pi_g \circ \text{NN}_\phi \circ \mathbb{R}_\psi \quad (16)$$

186 for all $g \in G$. Hence by Schur's lemma (Theorem 2), there exist a constant $C_{\phi, \psi} \in \mathbb{C}$ such that
 187 $\text{NN}_\phi \circ \mathbb{R}_\psi = C_{\phi, \psi} \text{Id}_{L^2(X)}$. Since $\text{NN}_\phi \circ \mathbb{R}_\psi$ is bilinear in ϕ and ψ , $C_{\phi, \psi}$ is bilinear in ϕ and ψ . \square

188 In particular, because depth- n feature map $\phi_{1:n}$ is G -equivariant (Lemma 2), the following depth- n
 189 \mathcal{H} -valued deep network $\text{DNN}[\gamma; \phi_{1:n}]$ is $L^2(X; \mathcal{H})$ -universal.

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

$$\text{DNN}[\gamma; \phi_{1:n}](x) := \int_{\Xi_1 \times \dots \times \Xi_n} \gamma(\xi_1, \dots, \xi_n) \phi_n(\bullet, \xi_n) \circ \dots \circ \phi_1(x, \xi_1) d\xi, \quad x \in X, \quad (17)$$

$$\mathbb{R}[f; \psi_{1:n}](\xi) := \int_{\Xi} \langle f(x), \psi_n(\bullet, \xi_n) \circ \dots \circ \psi_1(x, \xi_n) \rangle_{\mathcal{H}} dx, \quad \xi \in \Xi_1 \times \dots \times \Xi_n. \quad (18)$$

191 *Under the assumptions that $\text{DNN}_{\phi_{1:n}} \circ \mathbb{R}_{\psi_{1:n}}$ is bounded, and that π is irreducible, there exists a
 192 bilinear form $((\phi_{1:n}, \psi_{1:n}))$ satisfying $\text{DNN}_{\phi_{1:n}} \circ \mathbb{R}_{\psi_{1:n}} = ((\phi_{1:n}, \psi_{1:n})) \text{Id}_{L^2(X; \mathcal{H})}$.*

193 Again, it extends the original integral representation, and inherits the *linearization* trick of nonlinear
 194 parameters ξ by integrating all the possible parameters (beyond the difference of layers) and indirectly
 195 select which parameters to use by weighting on them.

196 4 Example: Depth- n Fully-Connected Network with Arbitrary Activation

197 As a concrete example, we present the ridgelet transform for depth- n fully-connected network.
 198 First, we show the depth-2 case based on a joint-affine-invariant argument, which was originally
 199 demonstrated by Sonoda et al. [33]. Then, we show the depth- n case based on a joint-equivariant
 200 argument by extending the original arguments.

201 We use the following known facts.

202 **Lemma 3.** *The regular representation π of the affine group $\text{Aff}(m)$ on $L^2(\mathbb{R}^m)$ (defined below) is
 203 irreducible.*

204 See Folland [10, Theorem 6.42] for the proof.

205 **Lemma 4.** *Suppose σ and ρ are a tempered distribution (S') and a Schwartz test function, respectively.
 206 Then, $S_\sigma \circ R_\rho : L^2(\mathbb{R}^m) \rightarrow L^2(\mathbb{R}^m)$ is bounded.*

207 See Sonoda et al. [29, Lemmas 7 and 8] for the proof.

208 **4.1 Depth-2**

209 Set $X := \mathbb{R}^m$ (data domain), $\Xi := \mathbb{R}^m \times \mathbb{R}$ (parameter domain), and $G := \text{Aff}(m) = GL(m) \ltimes \mathbb{R}^m$
 210 be the m -dimensional affine group, acting on data domain X by

$$g \cdot \mathbf{x} := L\mathbf{x} + \mathbf{t}, \quad g = (L, \mathbf{t}) \in GL(m) \ltimes \mathbb{R}^m, \quad \mathbf{x} \in X. \quad (19)$$

211 Addition to this, let π be the regular representation of $\text{Aff}(m)$ on $L^2(X)$, namely

$$\pi(g)[f](\mathbf{x}) := |\det L|^{-1/2} f(L^{-1}(\mathbf{x} - \mathbf{t})), \quad f \in L^2(X) \text{ and } g = (L, \mathbf{t}) \in GL(m) \ltimes \mathbb{R}^m. \quad (20)$$

212 Further, define a *dual action* of $\text{Aff}(m)$ on the parameter domain Ξ as

$$g \cdot (\mathbf{a}, b) = (L^{-\top} \mathbf{a}, b + \mathbf{t}^\top L^{-\top} \mathbf{a}), \quad g = (L, \mathbf{t}) \in GL(m) \ltimes \mathbb{R}^m, \quad (\mathbf{a}, b) \in \Xi. \quad (21)$$

213 Then, we can see the feature map $\phi(\mathbf{x}, (\mathbf{a}, b)) := \sigma(\mathbf{a} \cdot \mathbf{x} - b)$ is joint- G -invariant. In fact,

$$\phi(g \cdot \mathbf{x}, g \cdot (\mathbf{a}, b)) = \sigma(L^{-\top} \mathbf{a} \cdot (L\mathbf{x} + \mathbf{t}) - (b + \mathbf{t}^\top L^{-\top} \mathbf{a})) = \sigma(\mathbf{a} \cdot \mathbf{x} - b) = \phi(\mathbf{x}, (\mathbf{a}, b)).$$

214 By Lemma 3, the regular representation π of $G = \text{Aff}(m)$ is irreducible. Therefore, by Theorem 3,
 215 the depth-2 neural network and corresponding ridgelet transform:

$$\text{NN}[\gamma](\mathbf{x}) = \int_{\mathbb{R}^m \times \mathbb{R}} \gamma(\mathbf{a}, b) \sigma(\mathbf{a} \cdot \mathbf{x} - b) d\mathbf{a} db, \quad \text{and} \quad \mathbb{R}_2[f](\mathbf{a}, b) = \int_{\mathbb{R}^m} f(\mathbf{x}) \overline{\rho(\mathbf{a} \cdot \mathbf{x} - b)} d\mathbf{x},$$

216 satisfy the reconstruction formula $\text{NN} \circ \mathbb{R}_2 = ((\sigma, \rho)) \text{Id}_{L^2(\mathbb{R}^m)}$. In Appendix A, we supplemented a
 217 detailed proof. In Appendix B, we discussed a geometric interpretation of dual G -action (21).

218 **4.2 Depth- n**

219 Following Corollary 1, we derive the ridgelet transform for depth- n fully-connected network by
 220 constructing a joint-equivariant network.

221 Let $O(m)$ be the m -dimensional orthogonal group acting on \mathbb{R}^m by $Q\mathbf{v}$ for $Q \in O(m)$ and $\mathbf{v} \in \mathbb{R}^m$,
 222 and (re)set $G := O(m) \times \text{Aff}(m)$ be the product group, acting on the data domain X by

$$g \cdot \mathbf{x} := L\mathbf{x} + \mathbf{t}, \quad \mathbf{x} \in X, g = (Q, L, \mathbf{t}) \in G = O(m) \times (GL(m) \ltimes \mathbb{R}^m). \quad (22)$$

223 Namely, we set the $O(m)$ -action on X is trivial. Define a unitary representation π of G on vector-
 224 valued square-integrable functions $L^2(X; X)$ as

$$\pi_g[\mathbf{f}](\mathbf{x}) := Q\mathbf{f}(L^{-1}(\mathbf{x} - \mathbf{t})), \quad \mathbf{x} \in X, g = (Q, L, \mathbf{t}) \in G, \mathbf{f} \in L^2(X; X). \quad (23)$$

225 **Lemma 5.** *The above $\pi : G \rightarrow L^2(\mathbb{R}^m; \mathbb{R}^m)$ is an irreducible unitary representation.*

226 *Proof.* Recall that a tensor product of irreducible representations is irreducible. Since both $O(m)$ -
 227 action on \mathbb{R}^m and $\text{Aff}(m)$ -action on $L^2(\mathbb{R}^m)$ are irreducible, and $L^2(\mathbb{R}^m; \mathbb{R}^m)$ is a tensor product
 228 $\mathbb{R}^m \otimes L^2(\mathbb{R}^m)$, so the action π of product group $O(m) \times \text{Aff}(m)$ on tensor product $\mathbb{R}^m \otimes L^2(\mathbb{R}^m) =$
 229 $L^2(\mathbb{R}^m; \mathbb{R}^m)$ is irreducible. \square

230 Following the same arguments in Lemma 1, we first construct a *depth-2* joint- G -equivariant network.
 231 Take an arbitrary square-integrable (not yet joint- G -equivariant) vector-field $\mathbf{f}_0 \in L^2(X; X)$. Then,
 232 the following network is joint- G -equivariant:

$$\text{NN}(\mathbf{x}, \xi) := \text{NN}[\mathbb{R}_2[\pi_\xi[\mathbf{f}_0]]](\mathbf{x}) = \int_{\mathbb{R}^m \times \mathbb{R}} Q\mathbb{R}_2[\mathbf{f}_0](\mathbf{a}, b) \sigma(\mathbf{a}^\top L^{-1}(\mathbf{x} - \mathbf{t}) - b) d\mathbf{a} db, \quad (24)$$

233 for every $\mathbf{x} \in X, \xi = (Q, L, \mathbf{t}) \in O(m) \times GL(m) \ltimes \mathbb{R}^m$. Here \mathbb{R}_2 is the ridgelet transform for
 234 depth-2 case (applied for vector-valued function by element-wise manner). This is joint- G -equivariant
 235 because $\text{NN}(\mathbf{x}, \xi) = \pi_\xi[\mathbf{f}_0](\mathbf{x})$. Henceforth, we (re)set $\Xi := G$.

236 Finally, we construct a *depth- n* joint- G -equivariant network by composing the above depth-2 networks
 237 as below. Write $\xi := (\xi_1, \dots, \xi_n) \in \Xi^n$ for short. For any measurable function $\gamma : \Xi^n \rightarrow \mathbb{C}$ and
 238 vector-field $\mathbf{f} : \mathbb{R}^m \rightarrow \mathbb{R}^m$, put

$$\text{DNN}(\mathbf{x}) := \int_{\Xi^n} \gamma(\xi) \text{NN}(\bullet, \xi_n) \circ \dots \circ \text{NN}(\mathbf{x}, \xi_1) d\xi, \quad \mathbf{x} \in X \quad (25)$$

$$\mathbb{R}_n[\mathbf{f}](\xi) := \int_X \mathbf{f}(\mathbf{x})^\top \text{NN}(\bullet, \xi_n) \circ \dots \circ \text{NN}(\mathbf{x}, \xi_1) d\mathbf{x}, \quad \xi \in \Xi^n. \quad (26)$$

239 Then, as a consequence of Corollary 1, there exists a constant $c \in \mathbb{C}$ satisfying $\text{DNN} \circ \mathbb{R}_n[\mathbf{f}] = c\mathbf{f}$ for
 240 any $\mathbf{f} \in L^2(X; X)$.

241 **5 Example: Formal Deep Network**

242 We explain the *formal deep network* (FDN) introduced by Sonoda et al. [32]. Compared to the
 243 depth- n fully-connected network introduced in the previous section, the FDN (introduced in the
 244 previous study) is more abstract because the network architecture is not specified. Yet, we consider
 245 this is still useful for theoretical study of deep networks as it covers a wide range of groups and data
 246 domains (i.e., not limited to the affine group and the Euclidean space).

247 **5.1 Formal Deep Network**

248 Let G be an arbitrary locally compact group equipped with left-invariant measure dg , let X be a
 249 G -space equipped with left-invariant measure dx , and set $\Xi := G$ with right-invariant measure $d\xi$.
 250 The key concept is to identify each feature map $\phi : X \times \Xi \rightarrow X$ with a G -action $g : X \rightarrow X$ with
 251 parameter domain Ξ being identified with group G , and the composite of feature maps, say $g \circ h$,
 252 with product gh . Since a group is closed under its operation by definition, the proposed network can
 253 represent literally *any depth* such as a single hidden layer g , double hidden layers $g \circ h$, triple hidden
 254 layers $g \circ h \circ k$, and infinite hidden layers $g \circ h \circ \dots$. Besides, to lift the group action on a linear
 255 space, the network is formulated as a regular action of group G on a hidden layer, say $\psi \in L^2(X)$.

256 **Definition 6** (Formal Deep Network). For any functions $\psi \in L^2(X)$ and $\gamma : \Xi \rightarrow \mathbb{C}$, put

$$\text{DNN}[\gamma; \psi](x) := \int_{G_1 \times \dots \times G_n} \gamma(\xi_1, \dots, \xi_n) \psi \circ \xi_n \circ \dots \circ \xi_1(x) d\xi_1 \dots d\xi_n, \quad x \in X. \quad (27)$$

257 Here, $G = G_1 \times \dots \times G_n$ denotes the semi-direct product of groups, suggesting that the network
 258 gets much complex and expressive as it gets deeper.

259 To see the universality, define the dual action of G on the parameter domain $\Xi = G$ as

$$g \cdot \xi := \xi g^{-1}, \quad g \in G, \xi \in \Xi. \quad (28)$$

260 Then, we can see $\phi(x, \xi) := \psi \circ \xi(x)$ is joint- G -invariant. In fact,

$$\phi(g \cdot x, g \cdot \xi) = \psi \circ (g \cdot \xi)(g \cdot x) = \psi \circ (\xi \circ g^{-1})(g(x)) = \psi \circ \xi(x) = \phi(x, \xi).$$

261 Therefore, by Theorem 3, assuming that the regular representation $\pi : G \rightarrow \mathcal{U}(L^2(X))$ is irreducible,
 262 the ridgelet transform is given by

$$\mathbb{R}[f](\xi_1, \dots, \xi_n) = \int_X f(x) \overline{\psi \circ \xi_n \circ \dots \circ \xi_1(x)} dx, \quad (\xi_1, \dots, \xi_n) \in G_1 \times \dots \times G_n \quad (29)$$

263 satisfying $\text{NN} \circ \mathbb{R} = ((\sigma, \rho)) \text{Id}_{L^2(X)}$.

264 **5.2 Depth Separation**

265 To enjoy the advantage of abstract formulation, we discuss the effect of depth. For the sake of
 266 simplicity, we assume G to be a finite group, which may be acceptable given that the data domain
 267 X in practice is often discretized (or coarse-grained) into finite sets of representative points, say
 268 $X \approx \overline{X} := \{x_i\}_{i=1}^p$, and if so the G -action is also reduced to finite representative actions.

269 Following the concept of the formal deep network, we call group G acting on X a network. Let us
 270 consider depth-1 network G and depth- n network $G_1 \times \dots \times G_n$ satisfying $G = G_1 \times \dots \times G_n$. The
 271 equation indicates that two networks have the same expressive power, because they can implement
 272 the same class of maps $g : X \rightarrow X$.

273 Next, let us define the *width* of a single layer G as the cardinality $|G|$. This is reasonable because
 274 the set G parametrizes each map $g : X \rightarrow X$. Then, under the assumption that each G_i is simple,
 275 the depth- n network $G_1 \times \dots \times G_n$ can express the same class of depth-1 network exponentially-
 276 effectively, because the total widths are $\sum_{i=1}^n |G_i| = O(n)$ for depth- n and $\prod_{i=1}^n |G_i| = \exp O(n)$
 277 for depth-1. This estimate can be interpreted as the classical thought that the hierarchical models
 278 such as deep networks can represent complex functions combinatorially more efficient than shallow
 279 models.

280 6 Discussion

281 We have developed a systematic method for deriving a ridgelet transform for a wide range of learning
282 machines defined by joint-group-equivariant feature maps, yielding the universal approximation
283 theorems as corollaries. The previous results by Sonoda et al. [33] was limited to scalar-valued
284 joint-invariant functions, which were insufficient to deal with practical learning machines defined by
285 composite mappings of vector-valued functions, such as deep neural networks. For example, they
286 could only deal with abstract composite structures like formal deep network [32]. By extending their
287 argument to vector-valued joint-equivariant functions, we were able to deal with deep structures.
288 Traditionally, the techniques used in the expressive power analysis of deep networks were different
289 from those used in the analysis of shallow networks, as overviewed in the introduction. Nonetheless,
290 our main theorem cover both deep and shallow networks from the unified perspective (joint-group-
291 action on the data-parameter domain). Technically, this unification is due to Schur’s lemma, a basic
292 and useful result in the representation theory. Thanks to this lemma, the proof of the main theorem is
293 simple, yet the scope of application is wide. The significance of this study lies in revealing the close
294 relationship between machine learning theory and modern algebra. With this study as a catalyst, we
295 expect a major upgrade to machine learning theory from the perspective of modern algebra.

296 6.1 Limitations

297 In the main theorem, we assume the following: (1) joint-equivariance of feature map ϕ , (2) bound-
298 edness of composite operator $\mathbb{N} \circ \mathbb{R}$, (3) irreducibility of unitary representation π . In addition,
299 throughout this study, we assume (4) local compactness of group G , and (5) that the network is given
300 by the integral representation.

301 As discussed in the main text, satisfying (1) is much easier than (non-joint) equivariance. Also, (2) is
302 often a textbook exercise when the specific expression is given. (3) is required for Schur’s lemma, and
303 it is often sufficient to synthesize the known results such as the one for the example of depth- n fully-
304 connected network. (4) is quite a frequent assumption in the standard group representation theory, but
305 it excludes infinite-dimensional groups. When formulated *natively*, nonparametric learning models
306 including DNN can be infinite-dimensional groups. However, from the perspective of learnability,
307 it is nonsense to consider too large a model, and it is common to assume regularity conditions
308 such as sparsity and low rank in usual theoretical analysis. So, it is natural to impose additional
309 regularity conditions for satisfying local compactness. (5) may be rather an advantage because
310 there are established techniques to show the *cc*-universality of finite models by discretizing integral
311 representations. Moreover, there is a fast discretization scheme called the Barron’s rate based on the
312 quasi-Monte Carlo method. On the other hand, problems like the minimum width in the field of deep
313 narrow networks are analyses of finite parameters, and they could be a different type of parameters.
314 Yet, the current mainstream solutions are the information theoretic method by Park et al. [23] and the
315 neural ODE method by Cai [2], and both arguments contain the discretization of continuous models.
316 Therefore, we may expect a high affinity with the integral representation theory.

317 This study is the first step in extending the harmonic analysis method, which was previously applicable
318 only to shallow models, to deep models. The above limitations will be resolved in our future works.

319 7 Broader Impact

320 This work studies theoretical aspects of neural networks for expressing square integrable functions.
321 Since we do not propose a new method nor a new dataset, we expect that the impact of this work on
322 ethical aspects and future societal consequences will be small, if any. Our work can help understand
323 the theoretical benefit and limitations of neural networks in approximating functions. Our work and
324 the proof technique improve our understanding of the theoretical aspect of deep neural networks and
325 other learning machines used in machine learning, and may lead to better use of these techniques
326 with possible benefits to the society.

327 References

- 328 [1] M. M. Bronstein, J. Bruna, T. Cohen, and P. Veličković. Geometric Deep Learning: Grids, Groups, Graphs,
329 Geodesics, and Gauges. *arXiv preprint: 2104.13478*, 2021.

- 330 [2] Y. Cai. Achieve the Minimum Width of Neural Networks for Universal Approximation. In *The Eleventh*
331 *International Conference on Learning Representations*, 2023.
- 332 [3] E. J. Candès. *Ridgelets: theory and applications*. PhD thesis, Stanford University, 1998.
- 333 [4] R. T. Q. Chen, Y. Rubanova, J. Bettencourt, and D. Duvenaud. Neural Ordinary Differential Equations. In
334 *Advances in Neural Information Processing Systems*, volume 31, pages 6572–6583, Palais des Congrès de
335 Montréal, Montréal CANADA, 2018.
- 336 [5] L. Chizat and F. Bach. On the Global Convergence of Gradient Descent for Over-parameterized Models
337 using Optimal Transport. In *Advances in Neural Information Processing Systems 32*, pages 3036–3046,
338 Montreal, BC, 2018.
- 339 [6] A. Cohen, R. DeVore, G. Petrova, and P. Wojtaszczyk. Optimal Stable Nonlinear Approximation. *Founda-*
340 *tions of Computational Mathematics*, 22(3):607–648, 2022.
- 341 [7] N. Cohen, O. Sharir, and A. Shashua. On the Expressive Power of Deep Learning: A Tensor Analysis. In
342 *29th Annual Conference on Learning Theory*, volume 49, pages 1–31, 2016.
- 343 [8] I. Daubechies, R. DeVore, S. Foucart, B. Hanin, and G. Petrova. Nonlinear Approximation and (Deep)
344 ReLU Networks. *Constructive Approximation*, 55(1):127–172, 2022.
- 345 [9] W. E. A Proposal on Machine Learning via Dynamical Systems. *Communications in Mathematics and*
346 *Statistics*, 5(1):1–11, 2017.
- 347 [10] G. B. Folland. *A Course in Abstract Harmonic Analysis*. Chapman and Hall/CRC, New York, second
348 edition, 2015.
- 349 [11] P. Grohs, A. Klotz, and F. Voigtlaender. Phase Transitions in Rate Distortion Theory and Deep Learning.
350 *Foundations of Computational Mathematics*, 23(1):329–392, 2023.
- 351 [12] E. Haber and L. Ruthotto. Stable architectures for deep neural networks. *Inverse Problems*, 34(1):1–22,
352 2017.
- 353 [13] B. Hanin and M. Sellke. Approximating Continuous Functions by ReLU Nets of Minimal Width. *arXiv*
354 *preprint: 1710.11278*, 2017.
- 355 [14] P. Kidger and T. Lyons. Universal Approximation with Deep Narrow Networks. In *Proceedings of Thirty*
356 *Third Conference on Learning Theory*, volume 125 of *Proceedings of Machine Learning Research*, pages
357 2306–2327. PMLR, 2020.
- 358 [15] N. Kim, C. Min, and S. Park. Minimum width for universal approximation using ReLU networks on
359 compact domain. In *The Twelfth International Conference on Learning Representations*, 2024.
- 360 [16] L. Li, Y. Duan, G. Ji, and Y. Cai. Minimum Width of Leaky-ReLU Neural Networks for Uniform Universal
361 Approximation. In *Proceedings of the 40th International Conference on Machine Learning*, volume 202 of
362 *Proceedings of Machine Learning Research*, pages 19460–19470, 2023.
- 363 [17] Q. Li and S. Hao. An Optimal Control Approach to Deep Learning and Applications to Discrete-Weight
364 Neural Networks. In *Proceedings of The 35th International Conference on Machine Learning*, volume 80,
365 pages 2985–2994, Stockholm, 2018. PMLR.
- 366 [18] H. Lin and S. Jegelka. ResNet with one-neuron hidden layers is a Universal Approximator. In *Advances in*
367 *Neural Information Processing Systems*, volume 31, Montreal, BC, 2018.
- 368 [19] Z. Lu, H. Pu, F. Wang, Z. Hu, and L. Wang. The Expressive Power of Neural Networks: A View from the
369 Width. In *Advances in Neural Information Processing Systems*, volume 30, 2017.
- 370 [20] S. Mei, A. Montanari, and P.-M. Nguyen. A mean field view of the landscape of two-layer neural networks.
371 *Proceedings of the National Academy of Sciences*, 115(33):E7665–E7671, 2018.
- 372 [21] N. Murata. An integral representation of functions using three-layered networks and their approximation
373 bounds. *Neural Networks*, 9(6):947–956, 1996.
- 374 [22] A. Nitanda and T. Suzuki. Stochastic Particle Gradient Descent for Infinite Ensembles. *arXiv preprint:*
375 *1712.05438*, 2017.
- 376 [23] S. Park, C. Yun, J. Lee, and J. Shin. Minimum Width for Universal Approximation. In *International*
377 *Conference on Learning Representations*, 2021.

- 378 [24] G. Petrova and P. Wojtaszczyk. Limitations on approximation by deep and shallow neural networks.
379 *Journal of Machine Learning Research*, 24(353):1–38, 2023.
- 380 [25] G. Rotskoff and E. Vanden-Eijnden. Parameters as interacting particles: long time convergence and
381 asymptotic error scaling of neural networks. In *Advances in Neural Information Processing Systems 31*,
382 pages 7146–7155, Montreal, BC, 2018.
- 383 [26] J. W. Siegel. Optimal Approximation Rates for Deep ReLU Neural Networks on Sobolev and Besov
384 Spaces. *Journal of Machine Learning Research*, 24(357):1–52, 2023.
- 385 [27] S. Sonoda and N. Murata. Transportation analysis of denoising autoencoders: a novel method for analyzing
386 deep neural networks. In *NIPS 2017 Workshop on Optimal Transport & Machine Learning (OTML)*, pages
387 1–10, Long Beach, 2017.
- 388 [28] S. Sonoda, I. Ishikawa, and M. Ikeda. Ridge Regression with Over-Parametrized Two-Layer Networks
389 Converge to Ridgelet Spectrum. In *Proceedings of The 24th International Conference on Artificial
390 Intelligence and Statistics (AISTATS) 2021*, volume 130, pages 2674–2682. PMLR, 2021.
- 391 [29] S. Sonoda, I. Ishikawa, and M. Ikeda. Ghosts in Neural Networks: Existence, Structure and Role of
392 Infinite-Dimensional Null Space. *arXiv preprint: 2106.04770*, 2021.
- 393 [30] S. Sonoda, I. Ishikawa, and M. Ikeda. Universality of Group Convolutional Neural Networks Based
394 on Ridgelet Analysis on Groups. In *Advances in Neural Information Processing Systems 35*, pages
395 38680–38694, New Orleans, Louisiana, USA, 2022.
- 396 [31] S. Sonoda, I. Ishikawa, and M. Ikeda. Fully-Connected Network on Noncompact Symmetric Space
397 and Ridgelet Transform based on Helgason-Fourier Analysis. In *Proceedings of the 39th International
398 Conference on Machine Learning*, volume 162, pages 20405–20422, Baltimore, Maryland, USA, 2022.
- 399 [32] S. Sonoda, Y. Hashimoto, I. Ishikawa, and M. Ikeda. Deep Ridgelet Transform: Voice with Koopman
400 Operator Proves Universality of Formal Deep Networks. In *Proceedings of the 2nd NeurIPS Workshop on
401 Symmetry and Geometry in Neural Representations*, Proceedings of Machine Learning Research. PMLR,
402 2023.
- 403 [33] S. Sonoda, H. Ishi, I. Ishikawa, and M. Ikeda. Joint Group Invariant Functions on Data-Parameter Domain
404 Induce Universal Neural Networks. In *Proceedings of the 2nd NeurIPS Workshop on Symmetry and
405 Geometry in Neural Representations*, Proceedings of Machine Learning Research. PMLR, 2023.
- 406 [34] S. Sonoda, I. Ishikawa, and M. Ikeda. A unified Fourier slice method to derive ridgelet transform for a
407 variety of depth-2 neural networks. *Journal of Statistical Planning and Inference*, 233:106184, 2024.
- 408 [35] T. Suzuki. Generalization bound of globally optimal non-convex neural network training: Transportation
409 map estimation by infinite dimensional Langevin dynamics. In *Advances in Neural Information Processing
410 Systems 33*, pages 19224–19237, 2020.
- 411 [36] M. Telgarsky. Benefits of depth in neural networks. In *29th Annual Conference on Learning Theory*, pages
412 1–23, 2016.
- 413 [37] D. Yarotsky. Error bounds for approximations with deep ReLU networks. *Neural Networks*, 94:103–114,
414 2017.
- 415 [38] D. Yarotsky. Optimal approximation of continuous functions by very deep ReLU networks. In *Proceedings
416 of the 31st Conference On Learning Theory*, volume 75 of *Proceedings of Machine Learning Research*,
417 pages 639–649. PMLR, 2018.
- 418 [39] D. Yarotsky and A. Zhevnerchuk. The phase diagram of approximation rates for deep neural networks. In
419 *Advances in Neural Information Processing Systems*, volume 33, pages 13005–13015, 2020.

420 **A Depth-2 Fully-Connected Neural Network and Ridgelet Transform**

421 A non group theoretic proof by reducing to a Fourier expression is given in Sonoda et al. [29,
422 Theorem 6].

423 **A.1 Proof**

424 In the following, we identify the group G acting on data domain \mathbb{R}^m with the affine group $\text{Aff}(\mathbb{R}^m)$,
 425 and introduce the so-called twisted dual group action that leaves a function θ invariant. Then, we see
 426 the regular action π of G on functions space $L^2(\mathbb{R}^m)$ commutes with composite $S_\sigma \circ R_\rho$. Hence, by
 427 Schur's lemma, $S_\sigma \circ R_\rho$ is a constant multiple of identity, which concludes the assertion.

428 *Proof.* Let G be the affine group $\text{Aff}(\mathbb{R}^m) = GL(\mathbb{R}^m) \ltimes \mathbb{R}^m$. For any $g = (L, \mathbf{t}) \in G$, let

$$g \cdot \mathbf{x} := L\mathbf{x} + \mathbf{t}, \quad \mathbf{x} \in \mathbb{R}^m \quad (30)$$

429 be its action on \mathbb{R}^m , and let

$$\begin{aligned} \pi(g)[f](\mathbf{x}) &:= |\det L|^{-1/2} f(g^{-1} \cdot \mathbf{x}) \\ &= |\det L|^{-1/2} f(L^{-1}(\mathbf{x} - \mathbf{t})), \quad f \in L^2(\mathbb{R}^m) \end{aligned} \quad (31)$$

430 be its left-regular action on $L^2(\mathbb{R}^m)$.

431 Besides, putting

$$\theta((\mathbf{a}, b), \mathbf{x}) := \mathbf{a} \cdot \mathbf{x} - b, \quad (\mathbf{a}, b) \in \mathbb{R}^m \times \mathbb{R}, \mathbf{x} \in \mathbb{R}^m \quad (32)$$

432 we define the *twisted dual action* of G on $\mathbb{R}^m \times \mathbb{R}$ as

$$g \cdot (\mathbf{a}, b) := (L^{-\top} \mathbf{a}, b + \mathbf{a} \cdot (L^{-1} \mathbf{t})), \quad (\mathbf{a}, b) \in \mathbb{R}^m \times \mathbb{R} \quad (33)$$

433 so that the following invariance hold:

$$\theta(g \cdot (\mathbf{a}, b), g \cdot \mathbf{x}) = \theta((\mathbf{a}, b), \mathbf{x}) = \mathbf{a} \cdot \mathbf{x} - b. \quad (34)$$

434 To see this, use matrix expressions with extended variables

$$\theta((\mathbf{a}, b), \mathbf{x}) = (\mathbf{a}^\top \quad b) \begin{pmatrix} I_m & 0 \\ 0 & -1 \end{pmatrix} \begin{pmatrix} \mathbf{x} \\ 1 \end{pmatrix} =: \tilde{\mathbf{a}}^\top \tilde{I} \tilde{\mathbf{x}}, \quad (35)$$

$$g \cdot \tilde{\mathbf{x}} := \begin{pmatrix} g \cdot \mathbf{x} \\ 1 \end{pmatrix} = \begin{pmatrix} L & \mathbf{t} \\ 0 & 1 \end{pmatrix} \begin{pmatrix} \mathbf{x} \\ 1 \end{pmatrix} =: \tilde{L} \tilde{\mathbf{x}} \quad (36)$$

435 and calculate

$$\tilde{\mathbf{a}}^\top \tilde{I} \tilde{\mathbf{x}} = (\tilde{\mathbf{a}}^\top \tilde{I} \tilde{L}^{-1} \tilde{I}^{-1}) \tilde{I} (\tilde{L} \tilde{\mathbf{x}}) = (\tilde{I} \tilde{L}^{-\top} \tilde{I} \tilde{\mathbf{a}})^\top \tilde{I} (\tilde{L} \tilde{\mathbf{x}}), \quad (37)$$

436 which suggests $g \cdot \tilde{(\mathbf{a}, b)} := \tilde{I} \tilde{L}^{-\top} \tilde{I} \tilde{\mathbf{a}}$, and we have

$$\begin{aligned} \tilde{I} \tilde{L}^{-\top} \tilde{I} &= \begin{pmatrix} I_m & 0 \\ 0 & -1 \end{pmatrix} \begin{pmatrix} L & \mathbf{t} \\ 0 & 1 \end{pmatrix}^{-\top} \begin{pmatrix} I_m & 0 \\ 0 & -1 \end{pmatrix} \\ &= \begin{pmatrix} I_m & 0 \\ 0 & -1 \end{pmatrix} \begin{pmatrix} L^{-\top} & 0 \\ -\mathbf{t}^\top L^{-\top} & 1 \end{pmatrix} \begin{pmatrix} I_m & 0 \\ 0 & -1 \end{pmatrix} = \begin{pmatrix} L^{-\top} & 0 \\ \mathbf{t}^\top L^{-\top} & 1 \end{pmatrix}. \end{aligned}$$

437 Further, we define its regular-action by

$$\begin{aligned} \hat{\pi}(g)[\gamma](\mathbf{a}, b) &:= |\det L|^{1/2} \gamma(g^{-1} \cdot (\mathbf{a}, b)) \\ &= |\det L|^{1/2} \gamma(L^\top \mathbf{a}, b - \mathbf{a} \cdot \mathbf{t}), \quad (\mathbf{a}, b) \in \mathbb{R}^m \times \mathbb{R}. \end{aligned} \quad (38)$$

438 Then we can see that, for all $g = (L, \mathbf{t}) \in G$,

$$R_\rho \circ \pi(g) = \hat{\pi}(g) \circ R_\rho, \quad \text{and} \quad S_\sigma \circ \hat{\pi}(g) = \pi(g) \circ S_\sigma. \quad (39)$$

439 In fact, at every $g = (L, \mathbf{t}) \in G$ and $(\mathbf{a}, b) \in \mathbb{R}^m \times \mathbb{R}$,

$$R_\rho[\pi(g)[f]](\mathbf{a}, b) = |\det L|^{-1/2} \int_{\mathbb{R}^m} f(g^{-1} \cdot \mathbf{x}) \overline{\rho(\theta((\mathbf{a}, b), \mathbf{x}))} d\mathbf{x}$$

440 by putting $\mathbf{x} = g \cdot \mathbf{y} = L\mathbf{y} + \mathbf{t}$ with $d\mathbf{x} = |\det L| d\mathbf{y}$,

$$= |\det L|^{1/2} \int_{\mathbb{R}^m} f(\mathbf{y}) \overline{\rho(\theta((\mathbf{a}, b), g \cdot \mathbf{y}))} d\mathbf{y}$$

$$\begin{aligned}
&= |\det L|^{1/2} \int_{\mathbb{R}^m} f(\mathbf{y}) \overline{\rho(\theta(g^{-1} \cdot (\mathbf{a}, b), \mathbf{y})))} d\mathbf{y} \\
&= \widehat{\pi}(g)[R_\rho[f]](\mathbf{a}, b).
\end{aligned} \tag{40}$$

441 Similarly, at every $g = (L, \mathbf{t}) \in G$ and $\mathbf{x} \in \mathbb{R}^m$,

$$S_\sigma[\widehat{\pi}(g)[\gamma]](\mathbf{x}) = |\det L|^{1/2} \int_{\mathbb{R}^m \times \mathbb{R}} \gamma(g^{-1} \cdot (\mathbf{a}, b)) \sigma(\theta((\mathbf{a}, b), \mathbf{x})) d\mathbf{a}db$$

442 by putting $(\mathbf{a}, b) := g \cdot (\boldsymbol{\xi}, \eta) = (L^{-\top} \boldsymbol{\xi}, \eta + \boldsymbol{\xi} \cdot (L^{-1} \mathbf{t}))$ with $d\mathbf{a}db = |\det L| d\boldsymbol{\xi}d\eta$,

$$\begin{aligned}
&= |\det L|^{-1/2} \int_{\mathbb{R}^m \times \mathbb{R}} \gamma(\boldsymbol{\xi}, \eta) \sigma(\theta(g \cdot (\boldsymbol{\xi}, \eta), \mathbf{x})) d\boldsymbol{\xi}d\eta \\
&= |\det L|^{-1/2} \int_{\mathbb{R}^m \times \mathbb{R}} \gamma(\boldsymbol{\xi}, \eta) \sigma(\theta((\boldsymbol{\xi}, \eta), g^{-1} \cdot \mathbf{x})) d\boldsymbol{\xi}d\eta \\
&= \pi(g)[S_\sigma[\gamma]](\mathbf{x}).
\end{aligned} \tag{41}$$

443 Hence $S_\sigma \circ R_\rho$ commutes with $\pi(g)$ because

$$S_\sigma \circ R_\rho \circ \pi(g) = S_\sigma \circ \widehat{\pi}(g) \circ R_\rho = \pi(g) \circ S_\sigma \circ R_\rho.$$

444 Since $S_\sigma \circ R_\rho : L^2(\mathbb{R}^m) \rightarrow L^2(\mathbb{R}^m)$ is bounded (Lemma 4), and $(\pi, L^2(\mathbb{R}^m))$ is an irreducible
445 unitary representation of G (Lemma 3), Schur's lemma (Theorem 2) yields that there exist a constant
446 $C_{\sigma, \rho} \in \mathbb{C}$ such that

$$S_\sigma \circ R_\rho[f] = C_{\sigma, \rho} f \tag{42}$$

447 for all $f \in L^2(\mathbb{R}^m)$.

448 Finally, by directly computing the left-hand-side, namely $S_\sigma \circ R_\rho[f]$, we can verify that the constant
449 $C_{\sigma, \rho}$ is given by

$$C_{\sigma, \rho} = ((\sigma, \rho)) := (2\pi)^{m-1} \int_{\mathbb{R}} \sigma^\#(\omega) \overline{\rho^\#(\omega)} |\omega|^{-m} d\omega. \tag{43}$$

450

□

451 A.2 Proof for (33)

452 Use matrix expressions with extended variables

$$\theta((\mathbf{a}, b), \mathbf{x}) = (\mathbf{a}^\top \quad b) \begin{pmatrix} I_m & 0 \\ 0 & -1 \end{pmatrix} \begin{pmatrix} \mathbf{x} \\ 1 \end{pmatrix} =: \tilde{\mathbf{a}}^\top \tilde{I} \tilde{\mathbf{x}}, \tag{44}$$

$$\widetilde{g \cdot \mathbf{x}} := \begin{pmatrix} g \cdot \mathbf{x} \\ 1 \end{pmatrix} = \begin{pmatrix} L & \mathbf{t} \\ 0 & 1 \end{pmatrix} \begin{pmatrix} \mathbf{x} \\ 1 \end{pmatrix} =: \tilde{L} \tilde{\mathbf{x}} \tag{45}$$

453 and calculate

$$\tilde{\mathbf{a}}^\top \tilde{I} \tilde{\mathbf{x}} = (\tilde{\mathbf{a}}^\top \tilde{I} \tilde{L}^{-1} \tilde{I}^{-1}) \tilde{I}(\tilde{L} \tilde{\mathbf{x}}) = (\tilde{I} \tilde{L}^{-\top} \tilde{I} \tilde{\mathbf{a}})^\top \tilde{I}(\tilde{L} \tilde{\mathbf{x}}), \tag{46}$$

454 which suggests $g \cdot (\widetilde{\mathbf{a}}, b) := \tilde{I} \tilde{L}^{-\top} \tilde{I} \tilde{\mathbf{a}}$, and we have

$$\begin{aligned}
\tilde{I} \tilde{L}^{-\top} \tilde{I} &= \begin{pmatrix} I_m & 0 \\ 0 & -1 \end{pmatrix} \begin{pmatrix} L & \mathbf{t} \\ 0 & 1 \end{pmatrix}^{-\top} \begin{pmatrix} I_m & 0 \\ 0 & -1 \end{pmatrix} \\
&= \begin{pmatrix} I_m & 0 \\ 0 & -1 \end{pmatrix} \begin{pmatrix} L^{-\top} & 0 \\ -\mathbf{t}^\top L^{-\top} & 1 \end{pmatrix} \begin{pmatrix} I_m & 0 \\ 0 & -1 \end{pmatrix} = \begin{pmatrix} L^{-\top} & 0 \\ \mathbf{t}^\top L^{-\top} & 1 \end{pmatrix}.
\end{aligned}$$

455 **A.3 Proof for (39)**

456 In fact, at every $g = (L, \mathbf{t}) \in G$ and $(\mathbf{a}, b) \in \mathbb{R}^m \times \mathbb{R}$,

$$R_\rho[\pi(g)[f]](\mathbf{a}, b) = |\det L|^{-1/2} \int_{\mathbb{R}^m} f(g^{-1} \cdot \mathbf{x}) \overline{\rho(\theta((\mathbf{a}, b), \mathbf{x}))} d\mathbf{x}$$

457 by putting $\mathbf{x} = g \cdot \mathbf{y} = L\mathbf{y} + \mathbf{t}$ with $d\mathbf{x} = |\det L|d\mathbf{y}$,

$$\begin{aligned} &= |\det L|^{1/2} \int_{\mathbb{R}^m} f(\mathbf{y}) \overline{\rho(\theta((\mathbf{a}, b), g \cdot \mathbf{y}))} d\mathbf{y} \\ &= |\det L|^{1/2} \int_{\mathbb{R}^m} f(\mathbf{y}) \overline{\rho(\theta(g^{-1} \cdot (\mathbf{a}, b), \mathbf{y}))} d\mathbf{y} \\ &= \widehat{\pi}(g)[R_\rho[f]](\mathbf{a}, b). \end{aligned} \tag{47}$$

458 Similarly, at every $g = (L, \mathbf{t}) \in G$ and $\mathbf{x} \in \mathbb{R}^m$,

$$S_\sigma[\widehat{\pi}(g)[\gamma]](\mathbf{x}) = |\det L|^{1/2} \int_{\mathbb{R}^m \times \mathbb{R}} \gamma(g^{-1} \cdot (\mathbf{a}, b)) \sigma(\theta((\mathbf{a}, b), \mathbf{x})) d\mathbf{a}db$$

459 by putting $(\mathbf{a}, b) := g \cdot (\boldsymbol{\xi}, \eta) = (L^{-\top} \boldsymbol{\xi}, \eta + \boldsymbol{\xi} \cdot (L^{-1} \mathbf{t}))$ with $d\mathbf{a}db = |\det L|d\boldsymbol{\xi}d\eta$,

$$\begin{aligned} &= |\det L|^{-1/2} \int_{\mathbb{R}^m \times \mathbb{R}} \gamma(\boldsymbol{\xi}, \eta) \sigma(\theta(g \cdot (\boldsymbol{\xi}, \eta), \mathbf{x})) d\boldsymbol{\xi}d\eta \\ &= |\det L|^{-1/2} \int_{\mathbb{R}^m \times \mathbb{R}} \gamma(\boldsymbol{\xi}, \eta) \sigma(\theta((\boldsymbol{\xi}, \eta), g^{-1} \cdot \mathbf{x})) d\boldsymbol{\xi}d\eta \\ &= \pi(g)[S_\sigma[\gamma]](\mathbf{x}). \end{aligned} \tag{48}$$

460

461 **B Geometric Interpretation of Dual Action for Original Ridgelet Transform**

462 We explain a geometric interpretation of the dual action (33) in the previous section. We note that
 463 in general θ does not require any geometric interpretation as long as it is joint group invariant on
 464 data-parameter domain.

465 For each $(\mathbf{a}, b) \in \mathbb{R}^m \times \mathbb{R}$, put $\xi(\mathbf{a}, b) := \{\mathbf{x} \in \mathbb{R}^m \mid \mathbf{a} \cdot \mathbf{x} - b = 0\}$. Then it is a hyperplane in \mathbb{R}^m
 466 through point $\mathbf{x}_0 = b\mathbf{a}/|\mathbf{a}|^2$ with normal vector $\mathbf{u} := \mathbf{a}/|\mathbf{a}|$.

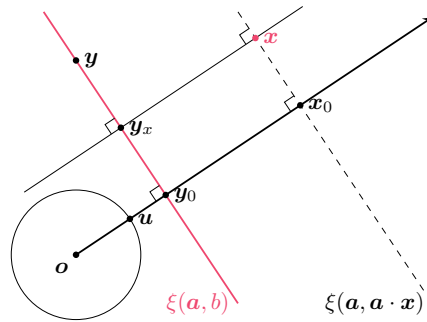


Figure 3: The invariant $\phi((\mathbf{a}, b), \mathbf{x}) = \sigma(\mathbf{a} \cdot \mathbf{x} - b)$ is the euclidean distance between point \mathbf{x} and hyperplane $\xi(\mathbf{a}, b)$ followed by scaling and nonlinearity σ

467 For any point \mathbf{y} in the hyperplane $\xi(\mathbf{a}, b)$, by definition $\mathbf{a} \cdot \mathbf{y} = b$, thus

$$\mathbf{a} \cdot \mathbf{x} - b = \mathbf{a} \cdot (\mathbf{x} - \mathbf{y}). \tag{49}$$

468 But this means $\mathbf{a} \cdot \mathbf{x} - b$ is a scaled distance between point \mathbf{x} and hyperplane $\xi(\mathbf{a}, b)$,

$$= |\mathbf{a}|d_E(\mathbf{x}, \xi(\mathbf{a}, b)), \tag{50}$$

and further a scaled distance between hyperplanes $\xi(\mathbf{a}, \mathbf{a} \cdot \mathbf{x})$ through \mathbf{x} with normal $\mathbf{a}/|\mathbf{a}|$ and $\xi(\mathbf{a}, b)$,

$$= |\mathbf{a}|d_E(\xi(\mathbf{a}, \mathbf{a} \cdot \mathbf{x}), \xi(\mathbf{a}, b)). \quad (51)$$

Now, we can interpret the invariant $\theta((\mathbf{a}, b), \mathbf{x}) := \mathbf{a} \cdot \mathbf{x} - b$ in a geometric manner, that is, it is the distance between point and hyperplane, or between hyperplanes. We note that we can regard entire $\sigma(\mathbf{a} \cdot \mathbf{x} - b)$ —the distance modulated by both scaling and nonlinearity—as the invariant, say ϕ .

Furthermore, the dual action $g \cdot (\mathbf{a}, b)$ is understood as a parallel translation of hyperplane $\xi(\mathbf{a}, b)$ to $\xi(g \cdot (\mathbf{a}, b))$ so as to leave the scaled distance θ invariant, namely

$$d_E(g \cdot \mathbf{x}, g \cdot \xi(\mathbf{a}, b)) = d_E(\mathbf{x}, \xi(\mathbf{a}, b)). \quad (52)$$

Indeed, for any $g = (L, \mathbf{t}) \in G$,

$$\begin{aligned} g \cdot \xi(\mathbf{a}, b) &= \{g \cdot \mathbf{x} \mid \mathbf{a} \cdot \mathbf{x} - b = 0\} \\ &= \{\mathbf{y} \mid \mathbf{a} \cdot (g^{-1} \cdot \mathbf{y}) - b = 0\} && \text{(by letting } \mathbf{y} = g \cdot \mathbf{x}\text{)} \\ &= \{\mathbf{y} \mid (L^{-\top}) \cdot \mathbf{y} - (b + \mathbf{a} \cdot (L^{-1}\mathbf{t})) = 0\} \\ &= \xi(g \cdot (\mathbf{a}, b)), \end{aligned}$$

meaning that the hyperplane with parameter (\mathbf{a}, b) translated by g is identical to the hyperplane with parameter $g \cdot (\mathbf{a}, b)$.

To summarize, in the case of fully-connected neural network (and its corresponding ridgelet transform), the invariant is a modulated distance $\sigma(\mathbf{a} \cdot \mathbf{x} - b)$, and the dual action is the parallel translation of hyperplane so as to keep the distance invariant. Further, from this geometric perspective, we can rewrite the fully-connected neural network in a geometric manner as

$$S[\gamma](\mathbf{x}) := \int_{\mathbb{R} \times \Xi} \gamma(\xi) \sigma(ad_E(\mathbf{x}, \xi)) da d\xi, \quad (53)$$

where $a \in \mathbb{R}$ denotes signed scale and Ξ denotes the space of all hyperplanes (not always through the origin). Since each hyperplane is parametrized by normal vectors $\mathbf{u} \in \mathbb{R}^{m-1}$ and distance $p \geq 0$ from the origin, we can induce the product of spherical measure $d\mathbf{u}$ and Lebesgue measure dp as a measure $d\xi$ on the space Ξ of hyperplanes.

486 **NeurIPS Paper Checklist**

487 **1. Claims**

488 Question: Do the main claims made in the abstract and introduction accurately reflect the
489 paper's contributions and scope?

490 Answer: [\[Yes\]](#)

491 Justification: Theorem 3 and Corollary 1

492 Guidelines:

- 493 • The answer NA means that the abstract and introduction do not include the claims
494 made in the paper.
- 495 • The abstract and/or introduction should clearly state the claims made, including the
496 contributions made in the paper and important assumptions and limitations. A No or
497 NA answer to this question will not be perceived well by the reviewers.
- 498 • The claims made should match theoretical and experimental results, and reflect how
499 much the results can be expected to generalize to other settings.
- 500 • It is fine to include aspirational goals as motivation as long as it is clear that these goals
501 are not attained by the paper.

502 **2. Limitations**

503 Question: Does the paper discuss the limitations of the work performed by the authors?

504 Answer: [\[Yes\]](#)

505 Justification: § 6.1

506 Guidelines:

- 507 • The answer NA means that the paper has no limitation while the answer No means that
508 the paper has limitations, but those are not discussed in the paper.
- 509 • The authors are encouraged to create a separate "Limitations" section in their paper.
- 510 • The paper should point out any strong assumptions and how robust the results are to
511 violations of these assumptions (e.g., independence assumptions, noiseless settings,
512 model well-specification, asymptotic approximations only holding locally). The authors
513 should reflect on how these assumptions might be violated in practice and what the
514 implications would be.
- 515 • The authors should reflect on the scope of the claims made, e.g., if the approach was
516 only tested on a few datasets or with a few runs. In general, empirical results often
517 depend on implicit assumptions, which should be articulated.
- 518 • The authors should reflect on the factors that influence the performance of the approach.
519 For example, a facial recognition algorithm may perform poorly when image resolution
520 is low or images are taken in low lighting. Or a speech-to-text system might not be
521 used reliably to provide closed captions for online lectures because it fails to handle
522 technical jargon.
- 523 • The authors should discuss the computational efficiency of the proposed algorithms
524 and how they scale with dataset size.
- 525 • If applicable, the authors should discuss possible limitations of their approach to
526 address problems of privacy and fairness.
- 527 • While the authors might fear that complete honesty about limitations might be used by
528 reviewers as grounds for rejection, a worse outcome might be that reviewers discover
529 limitations that aren't acknowledged in the paper. The authors should use their best
530 judgment and recognize that individual actions in favor of transparency play an impor-
531 tant role in developing norms that preserve the integrity of the community. Reviewers
532 will be specifically instructed to not penalize honesty concerning limitations.

533 **3. Theory Assumptions and Proofs**

534 Question: For each theoretical result, does the paper provide the full set of assumptions and
535 a complete (and correct) proof?

536 Answer: [\[Yes\]](#)

537 Justification: We put the proof right after Theorem 3

538 Guidelines:

- 539 • The answer NA means that the paper does not include theoretical results.
- 540 • All the theorems, formulas, and proofs in the paper should be numbered and cross-
- 541 referenced.
- 542 • All assumptions should be clearly stated or referenced in the statement of any theorems.
- 543 • The proofs can either appear in the main paper or the supplemental material, but if
- 544 they appear in the supplemental material, the authors are encouraged to provide a short
- 545 proof sketch to provide intuition.
- 546 • Inversely, any informal proof provided in the core of the paper should be complemented
- 547 by formal proofs provided in appendix or supplemental material.
- 548 • Theorems and Lemmas that the proof relies upon should be properly referenced.

549 4. Experimental Result Reproducibility

550 Question: Does the paper fully disclose all the information needed to reproduce the main ex-

551 perimental results of the paper to the extent that it affects the main claims and/or conclusions

552 of the paper (regardless of whether the code and data are provided or not)?

553 Answer: [NA]

554 Justification: This study does not include experiments.

555 Guidelines:

- 556 • The answer NA means that the paper does not include experiments.
- 557 • If the paper includes experiments, a No answer to this question will not be perceived
- 558 well by the reviewers: Making the paper reproducible is important, regardless of
- 559 whether the code and data are provided or not.
- 560 • If the contribution is a dataset and/or model, the authors should describe the steps taken
- 561 to make their results reproducible or verifiable.
- 562 • Depending on the contribution, reproducibility can be accomplished in various ways.
- 563 For example, if the contribution is a novel architecture, describing the architecture fully
- 564 might suffice, or if the contribution is a specific model and empirical evaluation, it may
- 565 be necessary to either make it possible for others to replicate the model with the same
- 566 dataset, or provide access to the model. In general, releasing code and data is often
- 567 one good way to accomplish this, but reproducibility can also be provided via detailed
- 568 instructions for how to replicate the results, access to a hosted model (e.g., in the case
- 569 of a large language model), releasing of a model checkpoint, or other means that are
- 570 appropriate to the research performed.
- 571 • While NeurIPS does not require releasing code, the conference does require all submis-
- 572 sions to provide some reasonable avenue for reproducibility, which may depend on the
- 573 nature of the contribution. For example
- 574 (a) If the contribution is primarily a new algorithm, the paper should make it clear how
- 575 to reproduce that algorithm.
- 576 (b) If the contribution is primarily a new model architecture, the paper should describe
- 577 the architecture clearly and fully.
- 578 (c) If the contribution is a new model (e.g., a large language model), then there should
- 579 either be a way to access this model for reproducing the results or a way to reproduce
- 580 the model (e.g., with an open-source dataset or instructions for how to construct
- 581 the dataset).
- 582 (d) We recognize that reproducibility may be tricky in some cases, in which case
- 583 authors are welcome to describe the particular way they provide for reproducibility.
- 584 In the case of closed-source models, it may be that access to the model is limited in
- 585 some way (e.g., to registered users), but it should be possible for other researchers
- 586 to have some path to reproducing or verifying the results.

587 5. Open access to data and code

588 Question: Does the paper provide open access to the data and code, with sufficient instruc-

589 tions to faithfully reproduce the main experimental results, as described in supplemental

590 material?

591
592
593
594
595
596
597
598
599
600
601
602
603
604
605
606
607
608
609
610
611
612
613
614
615
616
617
618
619
620
621
622
623
624
625
626
627
628
629
630
631
632
633
634
635
636
637
638
639
640
641
642

Answer: [NA] .

Justification: This study does not include experiments.

Guidelines:

- The answer NA means that paper does not include experiments requiring code.
- Please see the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- While we encourage the release of code and data, we understand that this might not be possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not including code, unless this is central to the contribution (e.g., for a new open-source benchmark).
- The instructions should contain the exact command and environment needed to run to reproduce the results. See the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- The authors should provide instructions on data access and preparation, including how to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- The authors should provide scripts to reproduce all experimental results for the new proposed method and baselines. If only a subset of experiments are reproducible, they should state which ones are omitted from the script and why.
- At submission time, to preserve anonymity, the authors should release anonymized versions (if applicable).
- Providing as much information as possible in supplemental material (appended to the paper) is recommended, but including URLs to data and code is permitted.

6. Experimental Setting/Details

Question: Does the paper specify all the training and test details (e.g., data splits, hyper-parameters, how they were chosen, type of optimizer, etc.) necessary to understand the results?

Answer: [NA]

Justification: This study does not include experiments

Guidelines:

- The answer NA means that the paper does not include experiments.
- The experimental setting should be presented in the core of the paper to a level of detail that is necessary to appreciate the results and make sense of them.
- The full details can be provided either with the code, in appendix, or as supplemental material.

7. Experiment Statistical Significance

Question: Does the paper report error bars suitably and correctly defined or other appropriate information about the statistical significance of the experiments?

Answer: [NA]

Justification: This study does not include experiments.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The authors should answer "Yes" if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.
- The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).
- The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)
- The assumptions made should be given (e.g., Normally distributed errors).
- It should be clear whether the error bar is the standard deviation or the standard error of the mean.

- 643
- 644
- 645
- 646
- 647
- 648
- 649
- 650
- It is OK to report 1-sigma error bars, but one should state it. The authors should preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of Normality of errors is not verified.
 - For asymmetric distributions, the authors should be careful not to show in tables or figures symmetric error bars that would yield results that are out of range (e.g. negative error rates).
 - If error bars are reported in tables or plots, The authors should explain in the text how they were calculated and reference the corresponding figures or tables in the text.

651 8. Experiments Compute Resources

652 Question: For each experiment, does the paper provide sufficient information on the com-
653 puter resources (type of compute workers, memory, time of execution) needed to reproduce
654 the experiments?

655 Answer: [NA]

656 Justification: This study does not include experiments.

657 Guidelines:

- 658
- 659
- 660
- 661
- 662
- 663
- 664
- 665
- The answer NA means that the paper does not include experiments.
 - The paper should indicate the type of compute workers CPU or GPU, internal cluster, or cloud provider, including relevant memory and storage.
 - The paper should provide the amount of compute required for each of the individual experimental runs as well as estimate the total compute.
 - The paper should disclose whether the full research project required more compute than the experiments reported in the paper (e.g., preliminary or failed experiments that didn't make it into the paper).

666 9. Code Of Ethics

667 Question: Does the research conducted in the paper conform, in every respect, with the
668 NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines?>

669 Answer: [Yes]

670 Justification: We have reviewed the NeurIPS Code of Ethics.

671 Guidelines:

- 672
- 673
- 674
- 675
- 676
- The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
 - If the authors answer No, they should explain the special circumstances that require a deviation from the Code of Ethics.
 - The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).

677 10. Broader Impacts

678 Question: Does the paper discuss both potential positive societal impacts and negative
679 societal impacts of the work performed?

680 Answer: [Yes]

681 Justification: § 7

682 Guidelines:

- 683
- 684
- 685
- 686
- 687
- 688
- 689
- 690
- 691
- 692
- 693
- The answer NA means that there is no societal impact of the work performed.
 - If the authors answer NA or No, they should explain why their work has no societal impact or why the paper does not address societal impact.
 - Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.
 - The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to

694 generate deepfakes for disinformation. On the other hand, it is not needed to point out
695 that a generic algorithm for optimizing neural networks could enable people to train
696 models that generate Deepfakes faster.

- 697 • The authors should consider possible harms that could arise when the technology is
698 being used as intended and functioning correctly, harms that could arise when the
699 technology is being used as intended but gives incorrect results, and harms following
700 from (intentional or unintentional) misuse of the technology.
- 701 • If there are negative societal impacts, the authors could also discuss possible mitigation
702 strategies (e.g., gated release of models, providing defenses in addition to attacks,
703 mechanisms for monitoring misuse, mechanisms to monitor how a system learns from
704 feedback over time, improving the efficiency and accessibility of ML).

705 11. Safeguards

706 Question: Does the paper describe safeguards that have been put in place for responsible
707 release of data or models that have a high risk for misuse (e.g., pretrained language models,
708 image generators, or scraped datasets)?

709 Answer: [NA]

710 Justification: This study does not contain any code, data nor trained model

711 Guidelines:

- 712 • The answer NA means that the paper poses no such risks.
- 713 • Released models that have a high risk for misuse or dual-use should be released with
714 necessary safeguards to allow for controlled use of the model, for example by requiring
715 that users adhere to usage guidelines or restrictions to access the model or implementing
716 safety filters.
- 717 • Datasets that have been scraped from the Internet could pose safety risks. The authors
718 should describe how they avoided releasing unsafe images.
- 719 • We recognize that providing effective safeguards is challenging, and many papers do
720 not require this, but we encourage authors to take this into account and make a best
721 faith effort.

722 12. Licenses for existing assets

723 Question: Are the creators or original owners of assets (e.g., code, data, models), used in
724 the paper, properly credited and are the license and terms of use explicitly mentioned and
725 properly respected?

726 Answer: [NA]

727 Justification: This study does not contain any code, data nor trained model

728 Guidelines:

- 729 • The answer NA means that the paper does not use existing assets.
- 730 • The authors should cite the original paper that produced the code package or dataset.
- 731 • The authors should state which version of the asset is used and, if possible, include a
732 URL.
- 733 • The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- 734 • For scraped data from a particular source (e.g., website), the copyright and terms of
735 service of that source should be provided.
- 736 • If assets are released, the license, copyright information, and terms of use in the
737 package should be provided. For popular datasets, paperswithcode.com/datasets
738 has curated licenses for some datasets. Their [licensing guide](#) can help determine the
739 license of a dataset.
- 740 • For existing datasets that are re-packaged, both the original license and the license of
741 the derived asset (if it has changed) should be provided.
- 742 • If this information is not available online, the authors are encouraged to reach out to
743 the asset's creators.

744 13. New Assets

745 Question: Are new assets introduced in the paper well documented and is the documentation
746 provided alongside the assets?

747
748
749
750
751
752
753
754
755
756
757
758
759
760
761
762
763
764
765
766
767
768
769
770
771
772
773
774
775
776
777
778
779
780
781
782
783
784
785
786
787
788
789
790
791

Answer: [NA]

Justification: This study does not provide any code, data nor trained model

Guidelines:

- The answer NA means that the paper does not release new assets.
- Researchers should communicate the details of the dataset/code/model as part of their submissions via structured templates. This includes details about training, license, limitations, etc.
- The paper should discuss whether and how consent was obtained from people whose asset is used.
- At submission time, remember to anonymize your assets (if applicable). You can either create an anonymized URL or include an anonymized zip file.

14. Crowdsourcing and Research with Human Subjects

Question: For crowdsourcing experiments and research with human subjects, does the paper include the full text of instructions given to participants and screenshots, if applicable, as well as details about compensation (if any)?

Answer: [NA]

Justification: This study does not involve crowdsourcing nor research with human subjects.

Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Including this information in the supplemental material is fine, but if the main contribution of the paper involves human subjects, then as much detail as possible should be included in the main paper.
- According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other labor should be paid at least the minimum wage in the country of the data collector.

15. Institutional Review Board (IRB) Approvals or Equivalent for Research with Human Subjects

Question: Does the paper describe potential risks incurred by study participants, whether such risks were disclosed to the subjects, and whether Institutional Review Board (IRB) approvals (or an equivalent approval/review based on the requirements of your country or institution) were obtained?

Answer: [NA]

Justification: This study does not involve crowdsourcing nor research with human subjects.

Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Depending on the country in which research is conducted, IRB approval (or equivalent) may be required for any human subjects research. If you obtained IRB approval, you should clearly state this in the paper.
- We recognize that the procedures for this may vary significantly between institutions and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for their institution.
- For initial submissions, do not include any information that would break anonymity (if applicable), such as the institution conducting the review.