
On Learning Sets of Symmetric Elements

Haggai Maron¹ Or Litany² Gal Chechik³ Ethan Fetaya³

Abstract

Learning from unordered sets is a fundamental learning setup, recently attracting increasing attention. Research in this area has focused on the case where elements of the set are represented by feature vectors, and far less emphasis has been given to the common case where set elements themselves adhere to their own symmetries. That case is relevant to numerous applications, from deblurring image bursts to multi-view 3D shape recognition and reconstruction. In this paper, we present a principled approach to learning sets of general symmetric elements. We first characterize the space of linear layers that are equivariant both to element reordering and to the inherent symmetries of elements, like translation in the case of images. We further show that networks that are composed of these layers, called *Deep Sets for Symmetric elements* layers (DSS), are universal approximators of both invariant and equivariant functions. DSS layers are also straightforward to implement. Finally, we show that they improve over existing set-learning architectures in a series of experiments with images, graphs and point-clouds.

1. Introduction

Learning with data that consists of unordered sets of elements is an important problem with numerous applications, from classification and segmentation of 3D data (Zaheer et al., 2017; Qi et al., 2017; Su et al., 2015; Kalogerakis et al., 2017) to image deblurring (Aittala & Durand, 2018). In this setting, each data point consists of a set of elements, and the task is independent of element order. This independence induces a symmetry structure, which can be used to design deep models with improved efficiency and generalization. Indeed, models that respect set symmetries,

e.g. (Zaheer et al., 2017; Qi et al., 2017), have become the leading approach for solving such tasks. However, in many cases, the elements of the set themselves adhere to certain symmetries, as happens when learning with sets of images, sets of point-clouds and sets of graphs. It is still unknown what is the best way to utilize these additional symmetries.

A common approach to handle per-element symmetries, is based on processing elements individually. First, one processes each set-element independently into a feature vector using a Siamese architecture (Bromley et al., 1994), and only then fuses information across all feature vectors. When following this process, the interaction between the elements of the set only occurs after each element has already been processed, possibly omitting low-level details. Indeed, it has been recently shown that for learning sets of images (Aittala & Durand, 2018; Sridhar et al., 2019; Liu et al., 2019), significant gain can be achieved with intermediate information-sharing layers.

In this paper, we present a principled approach to learning sets of symmetric elements. First, we describe the symmetry group of these sets, and then fully characterize the space of linear layers that are equivariant to this group. Notably, this characterization implies that information between set elements should be shared in all layers. For example, Figure 1 illustrates a DSS layer for sets of images. DSS layers provide a unified framework that generalizes several previously-described architectures for a variety of data types. In particular, it directly generalizes DeepSets (Zaheer et al., 2017). Moreover, other recent works can also be viewed as special cases of our approach (Hartford et al., 2018; Aittala & Durand, 2018; Sridhar et al., 2019).

A potential concern with equivariant architectures is that restricting layers to be equivariant to some group of symmetries may reduce the expressive power of the model (Maron et al., 2019c; Morris et al., 2018; Xu et al., 2019). We eliminate this potential limitation by proving two universal-approximation theorems for invariant and equivariant DSS networks. Simply put, these theorems state that if invariant (equivariant) networks for the elements of interest are universal, then the corresponding invariant (equivariant) DSS networks on sets of such elements are also universal.

To summarize, this paper has three main contributions: (1) We characterize the space of linear equivariant layers for sets

¹NVIDIA Research ²Stanford University ³Bar Ilan University.
Correspondence to: Haggai Maron <hmaron@nvidia.com>.

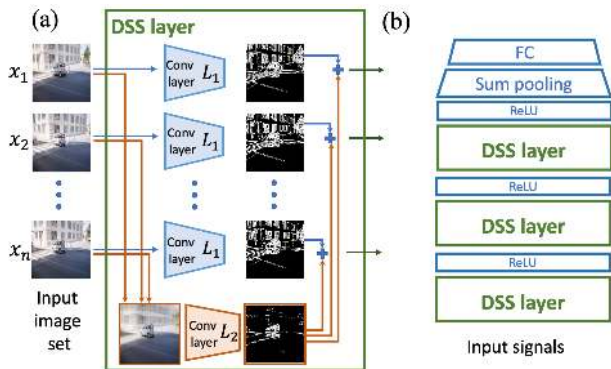


Figure 1. (a) A DSS layer for a set of images is composed of Siamese layer (blue) and an aggregation module (orange). The Siamese part is a convolutional layer (L_1) that is applied to each element independently. In the aggregation module, the *sum* of all images is processed by a different convolutional layer (L_2) and is added to the output of the Siamese part. (b) An example of a simple DSS-based invariant network.

of elements with symmetries. (2) We prove two universal approximation theorems for networks that are composed of DSS layers. (3) We demonstrate the empirical benefits of the DSS layers in a series of tasks, from classification through matching to selection, applied to diverse data from images to graphs and 3D point-clouds. These experiments show consistent improvement over previous approaches.

2. Previous work

Learning with sets. Several studies designed network architectures for set-structured input. Vinyals et al. (2015) suggested to extend the sequence-to-sequence framework of Sutskever et al. (2014) to handle sets. The prominent works of Ravanbakhsh et al. (2016); Edwards & Storkey (2016); Zaheer et al. (2017); Qi et al. (2017) proposed to use standard feed-forward neural networks whose layers are constrained to be equivariant to permutations. These models, when combined with a set-pooling layer, were also shown to be universal approximators of continuous permutation-invariant functions. Wagstaff et al. (2019) provided a theoretical study on the limitations of representing functions on sets with such networks. In another related work, Murphy et al. (2018) suggested to model permutation-invariant functions as an average of permutation-sensitive functions.

The specific case of learning sets of images was explored in several studies. Su et al. (2015); Kalogerakis et al. (2017) targeted classification and segmentation of 3D models by processing images rendered from several view points. These methods use a Siamese convolutional neural network to process the images, followed by view-pooling layer. Esteves et al. (2019) recently considered the same setup and sug-

gested to perform convolutions on a subgroup of the rotation group, which enables joint processing of all views. Sridhar et al. (2019) tackled 3D shape reconstruction from multiple view points and suggest using several equivariant mean-removal layers in which the mean of all images is subtracted from each image in the set. Aittala & Durand (2018) targeted image burst deblurring and denoising, and suggested to use set-pooling layers after convolutional blocks in which for each pixel, the maximum over all images is concatenated to all images. Liu et al. (2019) proposed to use an attention-based information sharing block for face recognition tasks. In Gordon et al. (2020) the authors modify neural processes by adding a translation equivariance assumption, treating the inputs as a set of translation equivariant objects.

Equivariance in deep learning. The prototypical example for equivariance in learning is probably visual object recognition, where the prevailing Convolutional Neural Networks (CNNs) are constructed from convolution layers which are equivariant to image translations. In the past few years, researchers have used invariance and equivariance considerations to devise deep learning architectures for other types of data. In addition to set-structured data discussed above, researchers suggested equivariant models for interaction between sets (Hartford et al., 2018), graphs (Kondor et al., 2018; Maron et al., 2019b;a; Chen et al., 2019; Albooyeh et al., 2019) and relational databases (Graham & Ravanbakhsh, 2019). Another successful line of work took into account other image symmetries such as reflections and rotations (Dieleman et al., 2016; Cohen & Welling, 2016a;b; Worrall et al., 2017; Cheng et al., 2018), spherical symmetries (Cohen et al., 2018; 2019b; Esteves et al., 2017), or 3D symmetries (Weiler et al., 2018; Winkels & Cohen, 2018; Worrall & Brostow, 2018; Kondor, 2018; Thomas et al., 2018; Weiler et al., 2018). From a theoretical point of view, several papers studied the properties of equivariant layers (Ravanbakhsh et al., 2017; Kondor & Trivedi, 2018; Cohen et al., 2019a) and characterized the expressive power of models that use such layers (Yarotsky, 2018; Maron et al., 2019c; Keriven & Peyré, 2019; Maehara & NT, 2019; Segol & Lipman, 2019).

3. Preliminaries

3.1. Notation and basic definitions

Let $x \in \mathbb{R}^\ell$ represent an input that adheres to a group of symmetries $G \leq S_\ell$, the symmetric group on ℓ elements. G captures those transformations that our task-of-interest is invariant (or equivariant) to. The action of G on \mathbb{R}^ℓ is defined by $(g \cdot x)_i = x_{g^{-1}(i)}$. For example, when inputs are images of size $h \times w$, we have $\ell = hw$ and G can be a group that applies cyclic translations, or left-right reflections to an image. A function is called G -equivariant if $f(g \cdot x) =$

$g \cdot f(x)$ for all $g \in G$. Similarly, a function f is called G -invariant if $f(g \cdot x) = f(x)$ for all $g \in G$.

3.2. G -invariant networks

G -equivariant networks are a popular way to model G -equivariant functions. These networks are composed of several linear G -equivariant layers, interleaved with activation functions like ReLU, and have the following form:

$$f = L_k \circ \sigma \circ L_{k-1} \cdots \circ \sigma \circ L_1, \quad (1)$$

Where $L_i : \mathbb{R}^{\ell \times d_i} \rightarrow \mathbb{R}^{\ell \times d_{i+1}}$ are linear G -equivariant layers, d_i are the feature dimensions and σ is a point-wise activation function. It is straightforward to show that this architecture results in a G -equivariant function. G -invariant networks are defined by adding an invariant layer on top of a G -equivariant function followed by a multilayer Perceptron (MLP), and have the form:

$$g = m \circ \sigma \circ h \circ \sigma \circ f, \quad (2)$$

where $h : \mathbb{R}^{\ell \times d_{k+1}} \rightarrow \mathbb{R}^{d_{k+2}}$ is a linear G -invariant layer and $m : \mathbb{R}^{d_{k+2}} \rightarrow \mathbb{R}^{d_{k+3}}$ is an MLP. It can be readily shown that this architecture results in a G -invariant function.

3.3. Characterizing equivariant layers

The main building block of G -invariant/equivariant networks are linear G -invariant/equivariant layers. To implement these networks, one has to characterize the space of linear G -invariant/equivariant layers, namely, L_i, h in Equations (1-2). For example, it is well known that for images with the group G of circular 2D translations, the space of linear G -equivariant layers is simply the space of all 2D convolutions operators (Puschel & Moura, 2008). Unfortunately, such elegant characterizations are not available for most permutation groups.

Characterizing linear G -equivariant layers can be reduced to the task of solving a set of linear equations in the following way: We are looking for a linear operator $L : \mathbb{R}^\ell \rightarrow \mathbb{R}^\ell$ that commutes with all the elements in G , namely:

$$L(g \cdot x) = g \cdot L(x), \quad x \in \mathbb{R}^\ell, \quad g \in G. \quad (3)$$

Note that L can be realized as a $\ell \times \ell$ matrix (which will be denoted in the same way), and as in Maron et al. (2019b), Equation 3 is equivalent to the following linear system:

$$g \cdot L = L, \quad g \in G, \quad (4)$$

where g acts on both dimensions of L . The solution space of Equation 4 characterizes the space of all G -equivariant linear layers, or equivalently, defines a parameter sharing scheme on the layer parameters for the group G (Wood & Shawe-Taylor, 1996; Ravanbakhsh et al., 2017). We will

denote the dimension of this space as $E(G)$. We note that in many important cases (e.g., (Zaheer et al., 2017; Hartford et al., 2018; Maron et al., 2019b; Albooyeh et al., 2019)) $|G|$ is exponential in ℓ so it is not possible to solve the linear system naively, and one has to resort to other strategies.

3.4. Deep Sets

Since the current paper generalizes *DeepSets* (Zaheer et al., 2017), we summarize their main results for completeness. Let $\{x_1, \dots, x_n\} \subset \mathbb{R}$ be a set, which we represent in arbitrary order as a vector $x \in \mathbb{R}^n$. DeepSets characterized all S_n -equivariant layers, namely, all matrices $L \in \mathbb{R}^{n \times n}$ such that $g \cdot L(x) = L(g \cdot x)$ for any permutation $g \in S_n$ and have shown that these operators have the following structure: $L = \lambda I_n + \beta \mathbf{1}\mathbf{1}^T$. When considering sets with higher dimensional features, i.e., $x_i \in \mathbb{R}^d$ and $X \in \mathbb{R}^{n \times d}$, this characterization takes the form:

$$L(X)_i = L_1(x_i) + L_2 \left(\sum_{j \neq i}^n x_j \right), \quad (5)$$

where $L_1, L_2 : \mathbb{R}^d \rightarrow \mathbb{R}^d$ are general linear functions and the subscript represents the i -th row of the output. The paper then suggests to concatenate several such layers, yielding a deep equivariant model (or an invariant model if a set pooling layer is added on top). Zaheer et al. (2017); Qi et al. (2017) established the universality of invariant networks that are composed of DeepSets Layers and Segol & Lipman (2019) extended this result to the equivariant case.

4. DSS layers

Our main goal is to design deep models for sets of elements with non-trivial per-element symmetries. In this section, we first formulate the symmetry group G of such sets. The deep models we advocate are composed of linear G -equivariant layers (DSS layers), therefore, our next step is to find a simple and practical characterization of the space of these layers.

4.1. Sets with symmetric elements

Let $\{x_1, \dots, x_n\} \subset \mathbb{R}^d$ be a set of elements with symmetry group $H \leq S_d$. We wish to characterize the space of linear maps $L : \mathbb{R}^{n \times d} \rightarrow \mathbb{R}^{n \times d}$ that are equivariant to both the natural symmetries of the elements, represented by the elements of the group H , as well as to the order of the n elements, represented by S_n .

In our setup, H operates on all elements x_i in the same way. More formally, the symmetry group is defined by $G = S_n \times H$, where S_n is the symmetric group on n elements. This group operates on $X \in \mathbb{R}^{n \times d}$ by applying the permutation $q \in S_n$ to

the first dimension and the same element $h \in H$ to the second dimension, namely $((q, h) \cdot X)_{ij} = X_{q^{-1}(i)h^{-1}(j)}$. Figure 2 illustrates this setup. Notably, this setup generalizes several popular learning setups: (1) DeepSets, where $H = \{I_d\}$ is the trivial group. (2) Tabular data (Hartford et al., 2018), where $H = S_d$. (3) Sets of images, where H is the group of circular translations (Aittala & Durand, 2018).

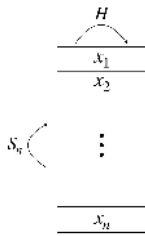


Figure 2. The input to a DSS layer is an $n \times d$ matrix, in which each row holds a d -dimensional element. $G = S_n \times H$ acts on it by applying a permutation to the columns and an element $h \in H$ to the rows.

One can also consider another setup, where the members of H that are applied to each element of the set may differ. Section C of the supplementary material formulates this setup and characterizes the corresponding equivariant layers in the common case where H acts transitively on $\{1, \dots, d\}$. While this setup can be used to model several interesting learning scenarios, it turns out that the corresponding equivariant networks are practically reduced to Siamese networks that were suggested in previous works.

4.2. Characterization of equivariant layers

This subsection provides a practical characterization of linear G -equivariant layers for $G = S_n \times H$. Our result generalizes DeepSets (equation 5) whose layers are tailored for $H = \{I_d\}$, by replacing the linear operators L_1, L_2 with linear H -equivariant operators. This result is summarized in the following theorem:

Theorem 1. Any linear G -equivariant layer $L : \mathbb{R}^{n \times d} \rightarrow \mathbb{R}^{n \times d}$ is of the form

$$L(X)_i = L_1^H(x_i) + L_2^H\left(\sum_{j \neq i}^n x_j\right),$$

where L_1^H, L_2^H are linear H -equivariant functions

Note that this is equivalent to the following formulations $L(X)_i = L_1^H(x_i) + L_2^H(\sum_{j=1}^n x_j) = L_1^H(x_i) + \sum_{j=1}^n L_2^H(x_j)$ due to linearity, and we will use them interchangeably throughout the paper. Figure 1 illustrates Theorem 1 for sets of images. In this case, applying a DSS layer amounts to: (i) Applying the same convolutional layer L_1 to all images in the set (blue); (ii) Applying another convolutional layer L_2 to the sum of all images (orange); and (iii) summing the outputs of these two layers. We discuss this theorem in the context of other widely-used data types such as point-clouds and graphs in section F of the Supplementary material.

We begin the proof by stating a useful lemma, that provides a formula for the dimension of the space of linear G -equivariant maps:

Lemma 1. Let $G \leq S_\ell$, then the dimension of the space of G -equivariant linear functions $L : \mathbb{R}^\ell \rightarrow \mathbb{R}^\ell$ is

$$E(G) = \frac{1}{|G|} \sum_{g \in G} \text{tr}(P(g))^2,$$

where $P(g)$ is the permutation matrix that corresponds to the permutation g .

The proof is given in the supplementary material. Given this lemma we can now prove Theorem 1:

Proof of Theorem 1. We wish to prove that all linear G -equivariant layers $L : \mathbb{R}^{n \times k} \rightarrow \mathbb{R}^{n \times k}$ are of the form $L(X)_i = L_1^H(x_i) + L_2^H(\sum_{j \neq i}^n x_j)$. Clearly, layers of this form are linear and equivariant. Moreover, the dimension of the space of these operators is exactly $2E(H)$ since we need to account for two linearly independent H -equivariant operators. The linear independence follows from the fact that their support in the matrix representation of L is disjoint. On the other hand, using Lemma 1 we have:

$$\begin{aligned} E(G) &= \frac{1}{|G|} \sum_{g \in G} \text{tr}(P(g))^2 = \\ &= \frac{1}{|H|} \frac{1}{n!} \sum_{q \in S_n} \sum_{h \in H} \text{tr}(P(q) \otimes P(h))^2 \\ &= \frac{1}{|H|} \frac{1}{n!} \sum_{q \in S_n} \sum_{h \in H} \text{tr}(P(q))^2 \text{tr}(P(h))^2 \\ &= \left(\frac{1}{|H|} \sum_{h \in H} \text{tr}(P(h))^2 \right) \cdot \left(\frac{1}{n!} \sum_{q \in S_n} \text{tr}(P(q))^2 \right) \\ &= E(H)E(S_n) = 2E(H). \end{aligned}$$

Here we used the fact that the trace is multiplicative with respect to the Kronecker product as well as the fact that $E(S_n) = 2$ (see (Zaheer et al., 2017) or Appendix 2 in (Maron et al., 2019b) for a generalization of this result).

To conclude, we have a linear subspace $\{L \mid L(X)_i = L_1^H(x_i) + L_2^H(\sum_{j \neq i}^n x_j)\}$, which is a subspace of the space of all linear G -equivariant operators, but has the same dimension, which implies that both spaces are equal. \square

Relation to (Aittala & Durand, 2018; Sridhar et al., 2019). In the specific case of a set of images and translation equivariance, L_i^H are convolutions. In this setting, (Aittala & Durand, 2018; Sridhar et al., 2019) have previously proposed using set-aggregation layers after convolutional

blocks. The main differences between these studies and the current paper are: (1) Our work applies to all types of symmetric elements and not just images; (2) We derive these layers from first principles; (3) We provide a theoretical analysis (Section 5); (4) We apply an aggregation step at each layer instead of only after convolutional blocks.

Generalizations. Section A of the supplementary material generalizes Theorem 1 to equivariant linear layers with multiple features. It also generalizes to several additional types of equivariant layers: $L : \mathbb{R}^{n \times d} \rightarrow \mathbb{R}$, $L : \mathbb{R}^{n \times d} \rightarrow \mathbb{R}^n$ and $L : \mathbb{R}^{n \times d} \rightarrow \mathbb{R}^d$. In addition, see Section B of the supplementary material for further discussion and characterization of the space of equivariant maps for a product of arbitrary permutation groups.

5. A universal approximation theorem

When restricting a network to be invariant (equivariant) to some group action, one may worry that these restrictions could reduce the network expressive power (see Maron et al. (2019c) or Xu et al. (2019) for concrete examples). We now show that networks that are constructed from DSS layers do not suffer from loss of expressivity. Specifically, we show that for any group H that induces a *universal* H -invariant (equivariant) network, its corresponding G -invariant (equivariant) network is universal as well.

We first state a lemma, which we later use for proving our universal-approximation theorems. The lemma shows that one can uniquely encode orbits of a group H in an invariant way by using a polynomial function. The full proof is given in Section D of the supplementary material.

Lemma 2. *Let $H \leq S_d$ then there exists a polynomial function $u : \mathbb{R}^d \rightarrow \mathbb{R}^l$, for some $l \in \mathbb{N}$, for which $u(x) = u(y)$ if and only if $x = h \cdot y$ for some $h \in H$.*

Proof idea. This lemma is a generalization of Proposition 1 in (Maron et al., 2019a) and we follow their proof. The main idea is that for any such group H there exists a finite set of invariant polynomials whose values on \mathbb{R}^d uniquely define each orbit of H in \mathbb{R}^d . \square

5.1. Invariant functions

We are now ready to state and prove our first universal approximation theorem. As before, the full proof can be found in the supplementary material (Section D).

Theorem 2. *Let $K \subset \mathbb{R}^{n \times d}$ be a compact domain such that $K = \cup_{g \in G} gK$. G -invariant networks are universal approximators (in $\|\cdot\|_\infty$ sense) of continuous G -invariant functions on K if and only if H -invariant networks are universal¹.*

¹We assume that there is a universal approximation theorem

Proof idea. The "only if" part is straightforward. For the "if" part, let $f : K \rightarrow \mathbb{R}$ be a continuous G -invariant function we wish to approximate. The idea of the proof is as follows: (1) we encode each element x_i with a unique H -invariant polynomial descriptor $u_H(x_i) \in \mathbb{R}^{l_H}$ (2) we encode the resulting set of descriptors with a unique S_n -invariant polynomial set descriptor $u_{S_n}(\{u_H(x_i)\}_{i \in [n]}) \in \mathbb{R}^{l_{S_n}}$ (3) we map the unique set descriptor $u_{S_n}(\{u(x_i)\}_{i \in [n]})$ to the appropriate value defined by f (4) we use the classic universal approximation theorem (Cybenko, 1989; Hornik et al., 1989) and our assumption on the universality of H -invariant networks to conclude that there exists a G -invariant network that can approximate each one of the previous stages to arbitrary precision on K . \square

Siamese networks. The proof of Theorem 1 implies that a simple Siamese architecture that applies an H -invariant network to each element in the set followed by a sum aggregation and finally an MLP is also universal. In section 6, we compare this architecture to our DSS networks and show that DSS-based architectures perform better in practice.

Relation to (Maron et al., 2019c). The authors proved that for any permutation group G , G -invariant networks have a universal approximation property, if the networks are allowed to use high-order tensors as intermediate representations (i.e., $X \in \mathbb{R}^d$ for $2 \leq l \leq n^2$), which are computationally prohibitive. We strengthen this result by proving that if first-order² H -invariant networks are universal, so are first-order G -invariant networks.

5.2. Equivariant functions

Three possible types of equivariant functions can be considered. First, functions of the form $f : \mathbb{R}^{n \times d} \rightarrow \mathbb{R}^n$. For example, such a function can model a selection task in which we are given a set $\{x_1, \dots, x_n\}$ and we wish to select a specific element from that set. Second, functions of the form $f : \mathbb{R}^{n \times d} \rightarrow \mathbb{R}^d$. An example for this type of functions would be an image-deblurring task in which we are given several noisy measurements of the same scene and we wish to generate a single high quality image (e.g., (Aittala & Durand, 2018)). Finally, functions of the form $f : \mathbb{R}^{n \times d} \rightarrow \mathbb{R}^{n \times d}$. This type of functions can be used to model tasks such as image co-segmentation where the input consists of several images and the task is to predict a joint segmentation map.

In this subsection we will prove a universality result for the third type of G -equivariant functions that were mentioned above, namely $f : \mathbb{R}^{n \times d} \rightarrow \mathbb{R}^{n \times d}$. We note that the equivariance of the first and second types can be easily

for the activation functions, e.g., ReLU.

²First-order networks use only first-order tensors.

deduced from this case. One can transform, for example, an $\mathbb{R}^{n \times d} \rightarrow \mathbb{R}^d$ G -equivariant function into a $\mathbb{R}^{n \times d} \rightarrow \mathbb{R}^{n \times d}$ function by repeating the \mathbb{R}^d vector n times and use our general approximation theorem on this function. We can get back a $\mathbb{R}^{n \times d} \rightarrow \mathbb{R}^d$ function by averaging over the first dimension.

Theorem 3. *Let $K \subset \mathbb{R}^{n \times d}$ be a compact domain such that $K = \cup_{g \in G} gK$. G -equivariant networks are universal approximators (in $\|\cdot\|_\infty$ sense) of continuous $\mathbb{R}^{n \times d} \rightarrow \mathbb{R}^{n \times d}$ G -equivariant functions on K if and only if H -equivariant networks are universal.*

Proof idea. The proof follows a similar line to the universality proof in (Segol & Lipman, 2019): First, we use the fact that equivariant polynomials are dense in the space of continuous equivariant functions. This enables us to assume that the function we wish to approximate is a G -equivariant polynomial. Next we show that for every output element, the mapping $\mathbb{R}^{n \times d} \rightarrow \mathbb{R}^d$ can be written as a sum of H -equivariant base polynomials with invariant coefficients. The base polynomials can be approximated by our assumption on H and the invariant mappings can be approximated by leveraging a slight modification of theorem 2. Finally we show how we can combine all the parts and approximate the full function with a G -equivariant network. \square

The full proof is given in Section D of the supplementary material. Similarly to the invariance case, using a Siamese network on each element separately followed by one DSS layer is sufficient for proving universality.

5.3. Examples

We can use Theorems (2-3) to show that DSS-based networks are universal in two important cases. For tabular data, which was considered by Hartford et al. (2018), the symmetries are $G = S_n \times S_d$. From the universality of S_n -invariant and equivariant networks (Zaheer et al., 2017; Segol & Lipman, 2019) we get that G -invariant (equivariant) networks are universal as well³. For sets of images, when H is the group of circular translations, it was shown in Yarotsky (2018) that H -invariant/equivariant networks are universal⁴, which implies universality of our DSS models.

6. Experiments

In this section we investigate the effectiveness of DSS layers in practice, by comparing them to previously suggested architectures and different aggregation schemes. We use the

³ Hartford et al. (2018) also considered interactions between more than two sets with $G = S_n \times S_{d_1} \times \dots \times S_{d_k}$. Our theorems can be extended to that case by induction on k .

⁴We note that this paper considers convolutional layers with full size kernels and no pooling layers

experiments to answer two basic questions: **(1) Early or late aggregation?** Can early aggregation architectures like DSS and its variants improve learning compared to Late aggregation architectures, which fuse the set information at the end of the data processing pipeline? and **(2) How to aggregate?** What is the preferred early aggregation scheme?

Tasks. We evaluated DSS in a series of six experiments spanning a wide range of tasks: from classification ($\mathbb{R}^{n \times d} \rightarrow \mathbb{R}$), through selection ($\mathbb{R}^{n \times d} \rightarrow \mathbb{R}^n$) and burst image deblurring ($\mathbb{R}^{n \times d} \rightarrow \mathbb{R}^d$) to general equivariant tasks ($\mathbb{R}^{n \times d} \rightarrow \mathbb{R}^{n \times d}$). The experiments also demonstrate the applicability of DSS to a range of data types, including point-clouds, images and graphs. Figure 3 illustrates the various types of tasks evaluated. A detailed description of all tasks, architectures and datasets is given in the supplementary material (Section E).

Competing methods. We compare DSS to four other models: (1) MLP; (2) DeepSets (DS) (Zaheer et al., 2017); (3) Siamese network; (4) Siamese network followed by DeepSets (Siamese+DS).

We also compare several variants of our DSS layers: **(1) DSS(sum):** our basic DSS layer from Theorem 1 **(2) DSS(max):** DSS with max-aggregation instead of sum-aggregation **(3) DSS(Aittala):** DSS with the aggregation proposed in (Aittala & Durand, 2018), namely, $L(x)_i \mapsto [L^H(x_i), \max_{j=1}^n L^H(x_j)]$ where $[]$ denotes feature concatenation and L^H is a linear H -equivariant layer **(4) DSS(Sridhar):** DSS layers with the aggregation proposed in (Sridhar et al., 2019), i.e., $L(x)_i \mapsto L^H(x_i) - \frac{1}{n} \sum_{j=1}^n L^H(x_j)$.

Evaluation protocol. For a fair comparison, for each particular task, all models have roughly the same number of parameters. In all experiments, we report the mean and standard deviation over 5 random initializations. Experiments were conducted using NVIDIA DGX with V100 GPUs.

6.1. Classification with multiple measurements

To illustrate the benefits of DSS, we first evaluate it in a signal-classification task using a synthetic dataset that we generated. Each sample consists of a set of $n = 25$ noisy measurements of the same 1D periodic signal sampled at 100 time-steps (see Figure 3). The clean signals are sampled uniformly from three signal types - sine, saw-tooth and square waves - with varying amplitude, DC component, phase-shift and frequency. The task is to predict the signal type given the set of noisy measurements. Figure 4 depicts the classification accuracy as a function of varying training set sizes, showing that DSS(sum) outperforms all other methods. Notably, DSS(sum) layers achieve significantly higher accuracy than the DeepSets architecture which

Dataset	Data type	Late Aggregation Siamese+DS	Early Aggregation				Random choice
			DSS (sum)	DSS (max)	DSS (Sridhar)	DSS (Aittala)	
UCF101	Images	36.41% ± 1.43	76.6% ± 1.51	76.39% ± 1.01	60.15% ± 0.76	77.96% ± 1.69	12.5%
Dynamic Faust	Point-clouds	22.26% ± 0.64	42.45% ± 1.32	28.71% ± 0.64	54.26% ± 1.66	26.43% ± 3.92	14.28%
Dynamic Faust	Graphs	26.53% ± 1.99	44.24% ± 1.28	30.54% ± 1.27	53.16% ± 1.47	26.66% ± 4.25	14.28%

Table 1. Frame selection tasks for images, point-clouds and graphs. Numbers represent average classification accuracy.

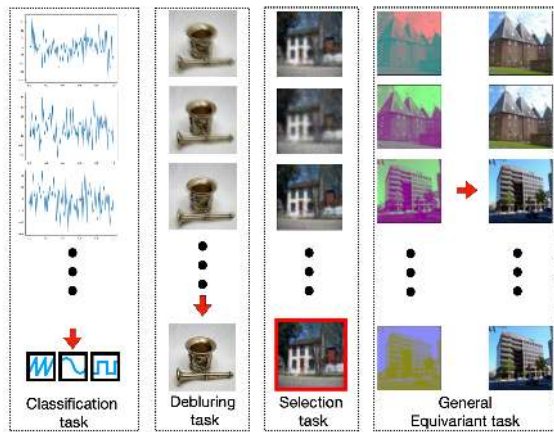


Figure 3. We consider all possible types of invariant and equivariant learning tasks in our settings: classification ($\mathbb{R}^{n \times d} \rightarrow \mathbb{R}$), selection ($\mathbb{R}^{n \times d} \rightarrow \mathbb{R}^n$), merging ($\mathbb{R}^{n \times d} \rightarrow \mathbb{R}^d$) and general equivariant tasks ($\mathbb{R}^{n \times d} \rightarrow \mathbb{R}^{n \times d}$).

takes into account the set structure but not within-element symmetry. DSS(sum) also outperforms the the *Siamese* and *Siamese+DS* architectures, which do not employ early aggregation. *DSS(Sridhar)* fails, presumably because it employs a mean removal aggregation scheme which is not appropriate for this task (removes the signal and leaves the noise).

6.2. Selection tasks

We next test DSS layers on selection tasks. In these tasks, we are given a set and wish to choose one element of the set that obeys a predefined property. Formally, each task is modelled as a G -equivariant function $f : \mathbb{R}^{n \times d} \rightarrow \mathbb{R}^n$, where the output vector represents the probability of selecting each element. The architecture comprises of three convolutional blocks employing Siamese or DSS variants, followed by a DeepSets block. We note that the *Siamese+DS* model was suggested for similar selection tasks in (Zaheer et al., 2017).

Frame selection in images and shapes. The first selection task is to find a particular frame within an unordered set of frames extracted from a video/shape sequence. For videos, we used the UCF101 dataset (Soomro et al., 2012). Each set contains $n = 8$ frames that were generated by randomly drawing a video, a starting position and frame ordering. The task is to select the "first" frame, namely, the one that appeared earliest in the video. Table 1 details the

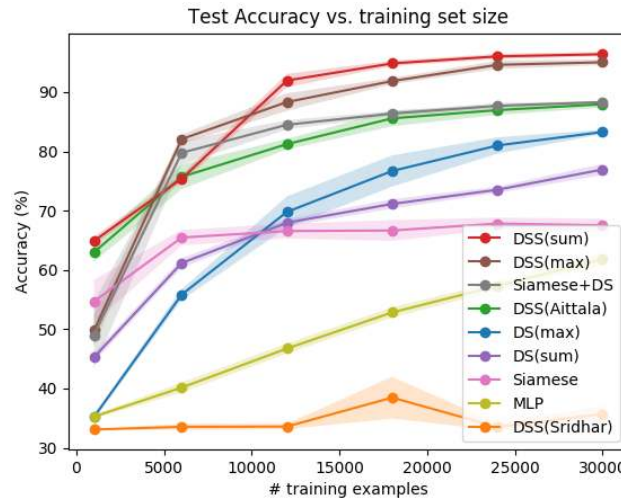


Figure 4. Comparison of set learning methods on the signal classification task. Shaded area represents standard deviation.

accuracy of all compared methods in this task, showing that *DSS(sum)* and *DSS(Aittala)* outperform *Siamese+DS* and *DSS(Sridhar)* by a large margin.

In a second selection task, we demonstrate that DSS can handle multiple data types. Specifically, we showcase how DSS operates on point-clouds and graphs. Given a short sequences of 3D human shapes performing various activities, the task is to identify which frame was the center frame in the original non-shuffled sequence. These human shapes are represented as point-clouds in the first experiment and as graphs (point-clouds + connectivity) in the second.

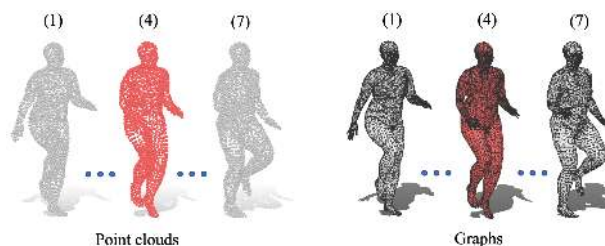


Figure 5. Shape-selection task on human shape sequences. Shapes are represented as graphs or as point-clouds. The task is to select the central frame (red). Numbers indicate frame order.

To generate the data, we cropped 7-frame-long sequences from the Dynamic Faust dataset (Bogo et al., 2017) in which

Noise type and strength	Late Aggregation Siamese+DS	Early Aggregation				Random choice
		DSS (sum)	DSS (max)	DSS (Sridahr)	DSS (Aittala)	
Gaussian $\sigma = 10$	77.2% \pm 0.37	78.48% \pm 0.48	77.99% \pm 1.1	76.8% \pm 0.25	78.34% \pm 0.49	5%
Gaussian $\sigma = 30$	65.89% \pm 0.66	68.35% \pm 0.55	67.85% \pm 0.40	61.52% \pm 0.54	66.89% \pm 0.58	5%
Gaussian $\sigma = 50$	59.24% \pm 0.51	62.6% \pm 0.45	61.59% \pm 1.00	55.25% \pm 0.40	62.02% \pm 1.03	5%
Occlusion 10%	82.15% \pm 0.45	83.13% \pm 1.00	83.27 \pm 0.51	83.21% \pm 0.338	83.19% \pm 0.67	5%
Occlusion 30%	77.47% \pm 0.37	78% \pm 0.89	78.69% \pm 0.32	78.71% \pm 0.26	78.27% \pm 0.67	5%
Occlusion 50%	76.2% \pm 0.82	77.29% \pm 0.40	76.64% \pm 0.45	77.04% \pm 0.75	77.03% \pm 0.58	5%

Table 2. Highest-quality image selection. Values indicate the mean accuracy.

Task	Late Aggregation Siamese+DS	Early Aggregation				TP
		DSS (sum)	DSS (max)	DSS (Sridahr)	DSS (Aittala)	
Color matching (places)	8.06 \pm 0.06	1.78 \pm 0.03	1.92 \pm 0.07	1.97 \pm 0.02	1.67 \pm 0.06	14.68
Color matching (CelebA)	6 \pm 0.13	1.27 \pm 0.07	1.34 \pm 0.07	1.35 \pm 0.03	1.17 \pm 0.04	18.72
Burst deblurring (Imagenet)	6.15 \pm 0.05	6.11 \pm 0.08	5.87 \pm 0.05	21.01 \pm 0.08	5.7 \pm 0.13	16.75

Table 3. Color-channel matching and burst deblurring tasks. Values indicate mean absolute error per pixel over the test set where the pixel values are in $[0, 255]$. TP stands for the trivial grey-scale predictor.

the shapes are given as triangular meshes. To generate point-clouds, we simply use the mesh vertices. To generate graphs, we use the graph defined by the triangular mesh⁵. See Figure 5 for an illustration of this task.

Results are summarized in Table 1, comparing DSS variants to a late-aggregation baseline (Siamese +DS) and to random choice. We further compared to a simple yet strong baseline. Using the mapping between points across shapes, we computed the mean of each point, and searched for the shape that was closest to that mean in L_1 sense. Frames in the sequence are $80msec$ apart, which limits the deviations around the mean, making it a strong baseline. Indeed, it achieved an accuracy of 34.47, which outperforms both late aggregation, DSS(max) and DSS(Aittala). In contrast, sum-based early aggregation methods reach significantly higher accuracy. Interestingly, using a graph representation provided a small improvement over point-clouds for almost all methods.

Highest quality image selection. Given a set of $n = 20$ degraded images of the same scene, the task is to select the highest-quality image. We generate data for this task from the Places dataset (Zhou et al., 2017), by adding noise and Gaussian blur to each image. The target image is defined to be the image that is the most similar in L_1 norm sense to the original image (see Figure 3 for an illustration). Notably, DSS consistently improves over Siamese+DS with a margin of 1% to 3%. See Table 2.

6.3. Color-channel matching

To illustrate the limitation of late-aggregation, we designed a very simple image-to-image task that highlights why early aggregation can be critical: learning to combine color chan-

nels into full images. Here, each sample consists of six images, generated from two randomly selected color images, by separating each image into three color channels. In each mono-chromatic image two channels were set to zero, yielding a $d = 64 \times 64 \times 3$ image. The task is to predict the fully colored image (i.e., imputing the missing color channels) for each of the set element. This can be formulated as a $\mathbb{R}^{n \times d} \rightarrow \mathbb{R}^{n \times d}$ G -equivariant task. See Figure 3 for an example.

We use a U-net architecture (Ronneberger et al., 2015), where convolutions and deconvolutions are replaced with Siamese layers or DSS variants. A DeepSets block is placed between the encoder and the decoder. Table 3 shows that layers with early aggregation significantly outperform DS+Siamese. For context, we add the error value of a trivial predictor which imputes the zeroed color channels by replicating the input color channel, resulting in a gray-scale image. This experiment was conducted on two datasets: *CelebA* (Liu et al., 2018), and *Places* (Zhou et al., 2017).

6.4. Burst image deblurring

Finally, we test DSS layers in a task of deblurring image bursts as in (Aittala & Durand, 2018). In this task, we are given a set of $n = 5$ blurred and noisy images of the same scene and aim to generate a single high quality image. This can be formulated as a $\mathbb{R}^{n \times d} \rightarrow \mathbb{R}^d$ G -equivariant task. See results in Table 3, where we also added the mean absolute error of a trivial predictor that outputs the median pixel of the images in the burst at each pixel. More details can be found in the supplementary material.

6.5. Summary of experiments

The above experiments demonstrate that applying early aggregation using DSS layers improves learning in various tasks and data types, compared with earlier architectures like *Siamese+DS*. More specifically, the basic DSS layer,

⁵In (Bogo et al., 2017) the points of each mesh are ordered consistently, providing point-to-point correspondence across frames. When this correspondence is not available, a shape matching algorithm like (Litany et al., 2017; Maron & Lipman, 2018) can be used as preprocessing.

$DSS(sum)$, performs well on all tasks, and $DSS(Aittala)$ has also yielded strong results. $DSS(Sridhar)$ performs well on some tasks but fails on others. See Section G of the supplementary materials for additional experiments on a multi-view reconstruction task.

Acknowledgments

This research was supported by an Israel science foundation grant 737/18. We thank Srinath Sridhar and Davis Remppe for useful discussions.

References

- Aittala, M. and Durand, F. Burst image deblurring using permutation invariant convolutional neural networks. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pp. 731–747, 2018.
- Albooyeh, M., Bertolini, D., and Ravanbakhsh, S. Incidence networks for geometric deep learning. *arXiv preprint arXiv:1905.11460*, 2019.
- Bogo, F., Romero, J., Pons-Moll, G., and Black, M. J. Dynamic FAUST: Registering human bodies in motion. In *IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, July 2017.
- Bromley, J., Guyon, I., LeCun, Y., Säckinger, E., and Shah, R. Signature verification using a” siamese” time delay neural network. In *Advances in neural information processing systems*, pp. 737–744, 1994.
- Chen, Z., Villar, S., Chen, L., and Bruna, J. On the equivalence between graph isomorphism testing and function approximation with gnns. *arXiv preprint arXiv:1905.12560*, 2019.
- Cheng, X., Qiu, Q., Calderbank, R., and Sapiro, G. Rotdcf: Decomposition of convolutional filters for rotation-equivariant deep networks. *arXiv preprint arXiv:1805.06846*, 2018.
- Cohen, T. and Welling, M. Group equivariant convolutional networks. In *International conference on machine learning*, pp. 2990–2999, 2016a.
- Cohen, T. S. and Welling, M. Steerable CNNs. (1990):1–14, 2016b. URL <http://arxiv.org/abs/1612.08498>.
- Cohen, T. S., Geiger, M., Köhler, J., and Welling, M. Spherical cnns. *arXiv preprint arXiv:1801.10130*, 2018.
- Cohen, T. S., Geiger, M., and Weiler, M. A general theory of equivariant cnns on homogeneous spaces. In *Advances in Neural Information Processing Systems*, pp. 9142–9153, 2019a.
- Cohen, T. S., Weiler, M., Kicanaoglu, B., and Welling, M. Gauge equivariant convolutional networks and the icosahedral cnn. *arXiv preprint arXiv:1902.04615*, 2019b.
- Cybenko, G. Approximation by superpositions of a sigmoidal function. *Mathematics of control, signals and systems*, 2(4):303–314, 1989.
- Deng, J., Dong, W., Socher, R., Li, L.-J., Li, K., and Fei-Fei, L. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pp. 248–255. Ieee, 2009.
- Dieleman, S., De Fauw, J., and Kavukcuoglu, K. Exploiting cyclic symmetry in convolutional neural networks. *arXiv preprint arXiv:1602.02660*, 2016.
- Edwards, H. and Storkey, A. Towards a neural statistician. *arXiv preprint arXiv:1606.02185*, 2016.
- Esteves, C., Allen-Blanchette, C., Makadia, A., and Daniilidis, K. 3d object classification and retrieval with spherical cnns. *arXiv preprint arXiv:1711.06721*, 2017.
- Esteves, C., Xu, Y., Allen-Blanchette, C., and Daniilidis, K. Equivariant multi-view networks. In *Proceedings of the IEEE International Conference on Computer Vision*, pp. 1568–1577, 2019.
- Fey, M. and Lenssen, J. E. Fast graph representation learning with pytorch geometric. *arXiv preprint arXiv:1903.02428*, 2019.
- Fulton, W. and Harris, J. *Representation theory: a first course*, volume 129. Springer Science & Business Media, 2013.
- Gordon, J., Bruinsma, W. P., Foong, A. Y. K., Requeima, J., Dubois, Y., and Turner, R. E. Convolutional conditional neural processes. In *International Conference on Learning Representations*, 2020. URL <https://openreview.net/forum?id=Skey4eBYP5>.
- Graham, D. and Ravanbakhsh, S. Deep models for relational databases. *arXiv preprint arXiv:1903.09033*, 2019.
- Hartford, J. S., Graham, D. R., Leyton-Brown, K., and Ravanbakhsh, S. Deep models of interactions across sets. In *ICML*, 2018.
- Hornik, K., Stinchcombe, M., and White, H. Multilayer feedforward networks are universal approximators. *Neural networks*, 2(5):359–366, 1989.
- Ioffe, S. and Szegedy, C. Batch normalization: Accelerating deep network training by reducing internal covariate shift. *arXiv preprint arXiv:1502.03167*, 2015.

- Kalogerakis, E., Averkiou, M., Maji, S., and Chaudhuri, S. 3d shape segmentation with projective convolutional networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 3779–3788, 2017.
- Keriven, N. and Peyré, G. Universal invariant and equivariant graph neural networks. *CoRR*, abs/1905.04943, 2019. URL <http://arxiv.org/abs/1905.04943>.
- Kingma, D. P. and Ba, J. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- Kipf, T. N. and Welling, M. Semi-supervised classification with graph convolutional networks. *arXiv preprint arXiv:1609.02907*, 2016.
- Kondor, R. N-body networks: a covariant hierarchical neural network architecture for learning atomic potentials. *arXiv preprint arXiv:1803.01588*, 2018.
- Kondor, R. and Trivedi, S. On the generalization of equivariance and convolution in neural networks to the action of compact groups. *arXiv preprint arXiv:1802.03690*, 2018.
- Kondor, R., Son, H. T., Pan, H., Anderson, B., and Trivedi, S. Covariant compositional networks for learning graphs. *arXiv preprint arXiv:1801.02144*, 2018.
- Litany, O., Remez, T., Rodolà, E., Bronstein, A., and Bronstein, M. Deep functional maps: Structured prediction for dense shape correspondence. In *Proceedings of the IEEE International Conference on Computer Vision*, pp. 5659–5667, 2017.
- Liu, X., Guo, Z., Li, S., Kong, L., Jia, P., You, J., and Kumar, B. Permutation-invariant feature restructuring for correlation-aware image set-based recognition. *arXiv preprint arXiv:1908.01174*, 2019.
- Liu, Z., Luo, P., Wang, X., and Tang, X. Large-scale celeb-faces attributes (celeba) dataset. Retrieved August, 15: 2018, 2018.
- Maehara, T. and NT, H. A simple proof of the universality of invariant/equivariant graph neural networks, 2019.
- Maron, H. and Lipman, Y. (probably) concave graph matching. In *Advances in Neural Information Processing Systems*, pp. 408–418, 2018.
- Maron, H., Ben-Hamu, H., Serviansky, H., and Lipman, Y. Provably powerful graph networks. *arXiv preprint arXiv:1905.11136*, 2019a.
- Maron, H., Ben-Hamu, H., Shamir, N., and Lipman, Y. Invariant and equivariant graph networks. In *International Conference on Learning Representations*, 2019b. URL <https://openreview.net/forum?id=Syx72jc9tm>.
- Maron, H., Fetaya, E., Segol, N., and Lipman, Y. On the universality of invariant networks. In *International conference on machine learning*, 2019c.
- Morris, C., Ritzert, M., Fey, M., Hamilton, W. L., Lenssen, J. E., Rattan, G., and Grohe, M. Weisfeiler and leman go neural: Higher-order graph neural networks. *arXiv preprint arXiv:1810.02244*, 2018.
- Murphy, R. L., Srinivasan, B., Rao, V., and Ribeiro, B. Janossy pooling: Learning deep permutation-invariant functions for variable-size inputs. *arXiv preprint arXiv:1811.01900*, 2018.
- Paszke, A., Gross, S., Chintala, S., Chanan, G., Yang, E., DeVito, Z., Lin, Z., Desmaison, A., Antiga, L., and Lerer, A. Automatic differentiation in pytorch. 2017.
- Puschel, M. and Moura, J. M. Algebraic signal processing theory: Foundation and 1-d time. *IEEE Transactions on Signal Processing*, 56(8):3572–3585, 2008.
- Qi, C. R., Su, H., Mo, K., and Guibas, L. J. Pointnet: Deep learning on point sets for 3d classification and segmentation. *Proc. Computer Vision and Pattern Recognition (CVPR)*, *IEEE*, 1(2):4, 2017.
- Ravanbakhsh, S., Schneider, J., and Póczos, B. Deep learning with sets and point clouds. *arXiv preprint arXiv:1611.04500*, 2016.
- Ravanbakhsh, S., Schneider, J., and Póczos, B. Equivariance through parameter-sharing. *arXiv preprint arXiv:1702.08389*, 2017.
- Ronneberger, O., Fischer, P., and Brox, T. U-net: Convolutional networks for biomedical image segmentation. In *International Conference on Medical image computing and computer-assisted intervention*, pp. 234–241. Springer, 2015.
- Segol, N. and Lipman, Y. On universal equivariant set networks. *arXiv preprint arXiv:1910.02421*, 2019.
- Simmons, G. F. *Introduction to topology and modern analysis*. Tokyo, 1963.
- Soomro, K., Zamir, A. R., and Shah, M. Ucf101: A dataset of 101 human actions classes from videos in the wild. *arXiv preprint arXiv:1212.0402*, 2012.
- Sridhar, S., Rempe, D., Valentin, J., Bouaziz, S., and Guibas, L. J. Multiview aggregation for learning category-specific shape reconstruction. *arXiv preprint arXiv:1907.01085*, 2019.
- Su, H., Maji, S., Kalogerakis, E., and Learned-Miller, E. Multi-view convolutional neural networks for 3d shape recognition. In *Proceedings of the IEEE international conference on computer vision*, pp. 945–953, 2015.

- Sutskever, I., Vinyals, O., and Le, Q. Sequence to sequence learning with neural networks. *Advances in NIPS*, 2014.
- Thomas, N., Smidt, T., Kearnes, S., Yang, L., Li, L., Kohlhoff, K., and Riley, P. Tensor field networks: Rotation-and translation-equivariant neural networks for 3d point clouds. *arXiv preprint arXiv:1802.08219*, 2018.
- Vinyals, O., Bengio, S., and Kudlur, M. Order matters: Sequence to sequence for sets. *arXiv preprint arXiv:1511.06391*, 2015.
- Wagstaff, E., Fuchs, F. B., Engelcke, M., Posner, I., and Osborne, M. On the limitations of representing functions on sets. *arXiv preprint arXiv:1901.09006*, 2019.
- Weiler, M., Geiger, M., Welling, M., Boomsma, W., and Cohen, T. 3D Steerable CNNs: Learning Rotationally Equivariant Features in Volumetric Data. 2018. URL <http://arxiv.org/abs/1807.02547>.
- Winkels, M. and Cohen, T. S. 3d g-cnns for pulmonary nodule detection. *arXiv preprint arXiv:1804.04656*, 2018.
- Wood, J. and Shawe-Taylor, J. Representation theory and invariant neural networks. *Discrete applied mathematics*, 69(1-2):33–60, 1996.
- Worrall, D. and Brostow, G. Cubenet: Equivariance to 3d rotation and translation. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pp. 567–584, 2018.
- Worrall, D. E., Garbin, S. J., Turmukhambetov, D., and Brostow, G. J. Harmonic networks: Deep translation and rotation equivariance. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 5028–5037, 2017.
- Xu, K., Hu, W., Leskovec, J., and Jegelka, S. How powerful are graph neural networks? In *International Conference on Learning Representations*, 2019. URL <https://openreview.net/forum?id=ryGs6iA5Km>.
- Yarotsky, D. Universal approximations of invariant maps by neural networks. *arXiv preprint arXiv:1804.10306*, 2018.
- Zaheer, M., Kottur, S., Ravanbakhsh, S., Póczos, B., Salakhutdinov, R. R., and Smola, A. J. Deep sets. In *Advances in neural information processing systems*, pp. 3391–3401, 2017.
- Zhou, B., Lapedriza, A., Khosla, A., Oliva, A., and Torralba, A. Places: A 10 million image database for scene recognition. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2017.