

---

# Set2Graph: Learning Graphs From Sets

---

Hadar Serviansky<sup>1</sup>

Nimrod Segol<sup>1</sup>

Jonathan Shlomi<sup>1</sup>

Kyle Cranmer<sup>2</sup>

Eilam Gross<sup>1</sup>

Haggai Maron<sup>3</sup>

Yaron Lipman<sup>1</sup>

<sup>1</sup>Weizmann Institute of Science <sup>2</sup>New York University <sup>3</sup>NVIDIA Research

## Abstract

Many problems in machine learning can be cast as learning functions from sets to graphs, or more generally to hypergraphs; in short, Set2Graph functions. Examples include clustering, learning vertex and edge features on graphs, and learning features on triplets in a collection.

A natural approach for building Set2Graph models is to characterize all linear equivariant set-to-hypergraph layers and stack them with non-linear activations. This poses two challenges: (i) the expressive power of these networks is not well understood; and (ii) these models would suffer from high, often intractable computational and memory complexity, as their dimension grows exponentially.

This paper advocates a family of neural network models for learning Set2Graph functions that is both practical and of maximal expressive power (universal), that is, can approximate arbitrary continuous Set2Graph functions over compact sets. Testing these models on different machine learning tasks, mainly an application to particle physics, we find them favorable to existing baselines.

## 1 Introduction

We consider the problem of learning functions taking sets of vectors in  $\mathbb{R}^{d_{in}}$  to graphs, or more generally hypergraphs; we name this problem Set2Graph, or set-to-graph. Set-to-graph functions appear in machine-learning applications such as clustering, predicting features on edges and nodes in graphs, and learning  $k$ -edge information in sets.

Mathematically, we represent each set-to-graph function as a collection of set-to- $k$ -edge functions, where each set-to- $k$ -edge function learns features on  $k$ -edges. That is, given an input set  $\mathcal{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_n\} \subset \mathbb{R}^{d_{in}}$  we consider functions  $\mathbf{F}^k$  attaching feature vectors to  $k$ -edges: each  $k$ -tuple  $(\mathbf{x}_{i_1}, \dots, \mathbf{x}_{i_k})$  is assigned with an output vector  $\mathbf{F}^k(\mathcal{X})_{i_1, i_2, \dots, i_k} \in \mathbb{R}^{d_{out}}$ . Now, functions mapping sets to hypergraphs with hyper-edges of size up-to  $k$  are modeled by  $(\mathbf{F}^1, \mathbf{F}^2, \dots, \mathbf{F}^k)$ . For example, functions mapping sets to standard graphs are represented by  $(\mathbf{F}^1, \mathbf{F}^2)$ , see Figure 1.

Set-to-graph functions are well-defined if they satisfy a property called *equivariance* (defined later), and therefore the set-to-graph problem is an instance of the bigger class of equivariant learning [3, 27, 16]. A natural approach for learning equivariant set-to-graph model is using out-of-the-box full equivariant model as in [20].

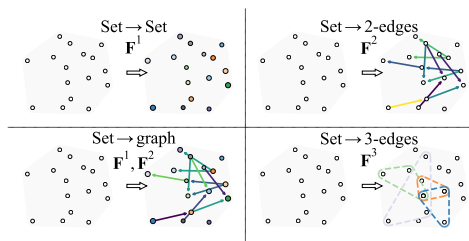


Figure 1: Set-to-graph functions are represented as collections of set-to- $k$ -edge functions.

A central question is: Are equivariant models *universal* for set-to-graph functions? That is, can equivariant models approximate any continuous equivariant function? In equivariant learning literature set-to-set models [44, 25] are proven equivariant universal [11, 30, 28]. In contrast, the situation for graph-to-graph equivariant models is more intricate: some models, such as message passing (a.k.a. graph convolutional networks), are known to be non-universal [41, 23, 19, 2], while high-order equivariant models are known to be universal [21] but require using high order tensors and therefore not practical. Universality of equivariant set-to-graph models is not known, as far as we are aware. In particular, are high order tensors required for universality (as the graph-to-graph case), or low order tensors (as in the set-to-set case) are sufficient?

In this paper we: (i) show that low order tensors are sufficient for set-to-graph universality, and (ii) build an equivariant model for the set-to-graph problem that is both *practical* (i.e., small number of parameters and no-need to build high-order tensors in memory) and *provably universal*. We achieve that with a composition of three networks:  $\mathbf{F}^k = \psi \circ \beta \circ \phi$ , where  $\phi$  is a set-to-set model,  $\beta$  is a *non-learnable* broadcasting set-to-graph layer, and  $\psi$  is a simple graph-to-graph network using only a single Multi-Layer Perceptron (MLP) acting on each  $k$ -edge independently.

Our main motivation for this work comes from an important set-to-2-edges learning problem in particle physics: partitioning (clustering) of simulated particles generated in the Large Hadron Collider (LHC). We demonstrate our model produces state of the art results on this task compared to relevant baselines. We also experimented with another set-to-2-edges problem of Delaunay triangulation, and a set-to-3-edges problem of 3D convex hull, in which we also achieve superior performances to the baselines.

## 2 Previous work

**Equivariant learning.** In many learning setups the task is invariant or equivariant to certain transformations of the input. The Canonical example is image recognition tasks [18, 17] and set classification tasks [44, 25]. Earlier methods such as [34] used non-equivariant methods to learn set functions. Restricting models to be invariant or equivariant to these transformation was shown to be an excellent approach for reducing the number of parameters of models while improving generalization [44, 25, 13, 9, 33, 41, 15, 20, 19, 3, 5, 7, 40, 4, 8, 37, 39, 37]. There has been a keen interest in the analysis of equivariant models [27, 15], especially the analysis of their approximation power [44, 25, 21, 11, 30, 22]. As far we know, the set-to-graph case was not treated before.

**Similarity learning.** Our work is related to the field of similarity learning, in which the goal is to learn a similarity function on pairs of inputs. In most cases, a siamese architecture is used in order to extract features for each input and then a similarity score is calculated based on this pair of features [1, 31, 43]. The difference from our setup is that similarity learning is aimed at extracting pairwise relations between two inputs, independently from the other members of the set, while we learn these pairwise relations globally from the entire input set. In the experimental section we show that the independence assumption taken in similarity learning might cause a significant degradation in performance compared to our global approach.

**Other related methods.** [10] suggest a method for meta-clustering that can be seen as an instance of the set2graph setup. Their method is based on LSTMs and therefore depends on the order of the set elements. In contrast, our method is blind (equivariant) to the chosen order of the input sets. The Neural relational Inference model [12] is another related work that targets learning relations and dynamics of a set of objects in an unsupervised manner. [35] had previously tackled planar Delaunay triangulation and convex hull prediction problems using a non-equivariant network.

## 3 Learning hypergraphs from sets

We would like to learn functions of sets of  $n$  vectors in  $\mathbb{R}^{d_{\text{in}}}$  to hypergraphs with  $n$  nodes (think of the nodes as corresponding to the set elements), and arbitrary  $k$ -edge feature vectors in  $\mathbb{R}^{d_{\text{out}}}$ , where a  $k$ -edge is defined as a  $k$ -tuple of set elements. A function mapping sets of vectors to  $k$ -edges is called set-to- $k$ -edge function and denoted  $\mathbf{F}^k : \mathbb{R}^{n \times d_{\text{in}}} \rightarrow \mathbb{R}^{n^k \times d_{\text{out}}}$ . Consequently, a set-to-hypergraph function would be modeled as a sequence  $(\mathbf{F}^1, \mathbf{F}^2, \dots, \mathbf{F}^K)$ , for target hypergraphs with hyperedges of maximal size  $K$ . For example,  $\mathbf{F}^2$  learns pairwise relations in a set; and  $(\mathbf{F}^1, \mathbf{F}^2)$  is a function from sets to graphs (outputs both node features and pairwise relations); see Figure 1.

**Our goal** is to design permutation equivariant neural network models for  $\mathbf{F}^k$  that are as-efficient-as-possible in terms of number of parameters and memory usage, but on the same time with maximal expressive power, i.e., universal.

**Representing sets and  $k$ -edges.** A matrix  $\mathbf{X} = (\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n)^T \in \mathbb{R}^{n \times d_{\text{in}}}$  represents a set of  $n$  vectors  $\mathbf{x}_i \in \mathbb{R}^{d_{\text{in}}}$  and therefore should be considered up to re-ordering of its rows. We denote by  $S_n = \{\sigma\}$  the symmetric group, that is the group of bijections (permutations)  $\sigma : [n] \rightarrow [n]$ , where  $[n] = \{1, \dots, n\}$ . We denote by  $\sigma \cdot \mathbf{X}$  the matrix resulting in reordering the rows of  $\mathbf{X}$  by the permutation  $\sigma$ , i.e.,  $(\sigma \cdot \mathbf{X})_{i,j} = \mathbf{X}_{\sigma^{-1}(i),j}$ . In this notation,  $\mathbf{X}$  and  $\sigma \cdot \mathbf{X}$  represent the same set, for all permutations  $\sigma$ .

$k$ -edges are represented as a tensor  $\mathbf{Y} \in \mathbb{R}^{n^k \times d_{\text{out}}}$ , where  $\mathbf{Y}_{i,:} \in \mathbb{R}^{d_{\text{out}}}$  denotes the feature vector attached to the  $k$ -edge defined by the  $k$ -tuple  $(\mathbf{x}_{i_1}, \mathbf{x}_{i_2}, \dots, \mathbf{x}_{i_k})$ , where  $i = (i_1, i_2, \dots, i_k) \in [n]^k$  is a multi-index with non-repeating indices. Similarly to the set case,  $k$ -edges are considered up-to renumbering of the nodes by some permutation  $\sigma \in S_n$ . That is, if we define the action  $\sigma \cdot \mathbf{Y}$  by  $(\sigma \cdot \mathbf{Y})_{i,j} = \mathbf{Y}_{\sigma^{-1}(i),j}$ , where  $\sigma^{-1}(i) = (\sigma^{-1}(i_1), \sigma^{-1}(i_2), \dots, \sigma^{-1}(i_k))$ , then  $\mathbf{Y}$  and  $\sigma \cdot \mathbf{Y}$  represent the same  $k$ -edge data, for all  $\sigma \in S_n$ .

**Equivariance.** A sufficient condition for  $\mathbf{F}^k$  to represent a well-defined map between sets  $\mathbf{X} \in \mathbb{R}^{n \times d_{\text{in}}}$  and  $k$ -edge data  $\mathbf{Y} \in \mathbb{R}^{n^k \times d_{\text{out}}}$  is *equivariance* to permutations, namely

$$\mathbf{F}^k(\sigma \cdot \mathbf{X}) = \sigma \cdot \mathbf{F}^k(\mathbf{X}), \quad (1)$$

for all sets  $\mathbf{X} \in \mathbb{R}^{n \times d_{\text{in}}}$  and permutations  $\sigma \in S_n$ . Equivariance guarantees, in particular, that the two equivalent sets  $\mathbf{X}$  and  $\sigma \cdot \mathbf{X}$  are mapped to equivalent  $k$ -edge data tensors  $\mathbf{F}^k(\mathbf{X})$  and  $\sigma \cdot \mathbf{F}^k(\mathbf{X})$ .

**Set-to- $k$ -edge models.** In this paper we explore the following equivariant neural network model family for approximating  $\mathbf{F}^k$ :

$$\mathbf{F}^k(\mathbf{X}; \theta) = \psi \circ \beta \circ \phi(\mathbf{X}), \quad (2)$$

where  $\phi, \beta$ , and  $\psi$  will be defined soon. For  $\mathbf{F}^k$  to be equivariant (as in equation 1) it is sufficient that its constituents, namely  $\phi, \beta, \psi$ , are equivariant. That is,  $\phi, \beta, \psi$  all satisfy equation 1.

**Set-to-graphs models.** Given the model of set-to- $k$ -edge functions, a model for a set-to-graph function can now be constructed from a pair of set-to- $k$ -edge networks  $(\mathbf{F}^1, \mathbf{F}^2)$ . Similarly, set-to-hypergraph function would require  $(\mathbf{F}^1, \dots, \mathbf{F}^K)$ , where  $K$  is the maximal hyperedge size. Figure 1 shows an illustration of set-to- $k$ -edge and set-to-graph functions

**$\phi$  component.**  $\phi : \mathbb{R}^{n \times d_{\text{in}}} \rightarrow \mathbb{R}^{n \times d_1}$  is a set-to-set equivariant model, that is  $\phi$  is mapping sets of vectors in  $\mathbb{R}^{d_{\text{in}}}$  to sets of vectors in  $\mathbb{R}^{d_1}$ . To achieve the universality goal we will need  $\phi$  to be universal as set-to-set model; that is,  $\phi$  can approximate arbitrary continuous set-to-set functions. Several options exists [11, 28] although probably the simplest option is either DeepSets [44] or one of its variations; all were proven to be universal recently in [30].

In practice, as will be clear later from the proof of the universality of the model, when building set-to-graph or set-to-hypergraph model, the  $\phi$  (set-to-set) part of the  $k$ -edge networks can be shared between different set-to- $k$ -edge models,  $\mathbf{F}^k$ , without compromising universality.

**$\beta$  component.**  $\beta : \mathbb{R}^{n \times d_1} \rightarrow \mathbb{R}^{n^k \times d_2}$  is a non-learnable linear *broadcasting layer* mapping sets to  $k$ -edges. In theory, as shown in [20] the space of equivariant linear mappings  $\mathbb{R}^{n \times d_1} \rightarrow \mathbb{R}^{n^k \times d_2}$  is of dimension  $d_1 d_2 \text{bell}(k+1)$  which can be very high since bell numbers have exponential growth. Interestingly, in the set-to- $k$ -edge case one can achieve universality with only  $k$  linear operators. We define the broadcasting operator to be

$$\beta(\mathbf{X})_{i,:} = [\mathbf{x}_{i_1}, \mathbf{x}_{i_2}, \dots, \mathbf{x}_{i_k}], \quad (3)$$

where  $i = (i_1, \dots, i_k)$  and brackets denote concatenation in the feature dimension, that is, for  $\mathbf{A} \in \mathbb{R}^{n^k \times d_a}$ ,  $\mathbf{B} \in \mathbb{R}^{n^k \times d_b}$  their concatenation is  $[\mathbf{A}, \mathbf{B}] \in \mathbb{R}^{n^k \times (d_a + d_b)}$ . Therefore, the feature output dimension of  $\beta$  is  $d_2 = k d_1$ .

As an example, consider the graph case, where  $k = 2$ . In this case  $\beta(\mathbf{X})_{i_1, i_2, :} = [\mathbf{x}_{i_1}, \mathbf{x}_{i_2}]$ . This function is illustrated in Figure 2 broadcasting data in  $\mathbb{R}^{n \times d_1}$  to tensor  $\mathbb{R}^{n \times n \times d_2}$ .

To see that the broadcasting layer is equivariant, it is enough to consider a single feature  $\beta(\mathbf{X})_i = \mathbf{x}_{i_1}$ . Permuting the rows of  $\mathbf{X}$  by a permutation  $\sigma$  we get  $\beta(\sigma \cdot \mathbf{X})_{i, j} = \mathbf{x}_{\sigma^{-1}(i_1), j} = \beta(\mathbf{X})_{\sigma^{-1}(i), j} = (\sigma \cdot \beta(\mathbf{X}))_{i, j}$ .

**$\psi$  component.**  $\psi : \mathbb{R}^{n^k \times d_2} \rightarrow \mathbb{R}^{n^k \times d_{\text{out}}}$  is a mapping of  $k$ -tensors to  $k$ -tensors. Here the theory of equivariant operators indicates that the space of linear equivariant maps is of dimension  $d_2 d_{\text{out}} \text{bell}(2k)$  that suggests a huge number of model parameters even for a single linear layer. Surprisingly, universality can be achieved with much less, in fact a single linear operator (i.e., scaled identity) in each layer. In the multi-feature multi-layer case this boils to applying a Multi-Layer Perceptron  $\mathbf{m} : \mathbb{R}^{d_2} \rightarrow \mathbb{R}^{d_{\text{out}}}$  independently to each feature vector in the input tensor  $\mathbf{X} \in \mathbb{R}^{n^k \times d_2}$ . That is, we use

$$\psi(\mathbf{X})_{i, :} = \mathbf{m}(\mathbf{X}_{i, :}). \quad (4)$$

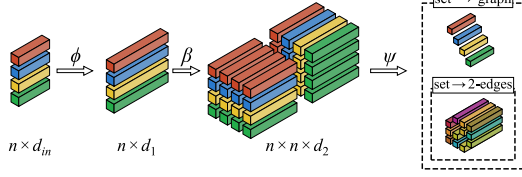


Figure 2: The model architecture for the Set-to-graph and set-to-2-edge functions.

Figure 2 illustrates set-to-2-edges and set-to-graph models incorporating the three components  $\phi, \beta, \psi$  discussed above. We note that, Indeed,  $\phi, \beta, \psi$  are equivariant.

## 4 Universality of set-to-graph models

In this section we prove that the model  $\mathbf{F}^k$  introduced above, is universal, in the sense it can approximate arbitrary continuous equivariant set-to- $k$ -edge functions  $\mathbf{G}^k : \mathbb{R}^{n \times d_{\text{in}}} \rightarrow \mathbb{R}^{n^k \times d_{\text{out}}}$  over compact domains  $K \subset \mathbb{R}^{n \times d_{\text{in}}}$ .

**Theorem 1.** *The model  $\mathbf{F}^k$  is set-to- $k$ -edge universal.*

A corollary of Theorem 1 establishes set-to-hypergraph universal models:

**Theorem 2.** *The model  $(\mathbf{F}^1, \dots, \mathbf{F}^k)$  is set-to-hypergraph universal.*

Our main tool for proving Theorem 1 is a characterization of the equivariant set-to- $k$ -edge polynomials  $\mathbf{P}^k$ . This characterization can be seen as a generalization of the characterization of set-to-set equivariant polynomial recently appeared in [30].

We consider an arbitrary set-to- $k$ -edge continuous mapping  $\mathbf{G}^k(\mathbf{X})$  over a compact set  $K \subset \mathbb{R}^{n \times d_{\text{in}}}$ . Since  $\mathbf{G}^k$  is equivariant we can assume  $K$  is symmetric, i.e.,  $\sigma \cdot K = K$  for all  $\sigma \in S_n$ . The proof consists of three parts: (i) Characterization of the equivariant set-to- $k$ -edge polynomials  $\mathbf{P}^k$ . (ii) Showing that every equivariant continuous set-to- $k$ -edge function  $\mathbf{G}^k$  can be approximated by some  $\mathbf{P}^k$ . (iii) Every  $\mathbf{P}^k$  can be approximated by our model  $\mathbf{F}^k$ .

Before providing the full proof which contains some technical derivations let us provide a simpler universality proof (under some mild extra conditions) for the set-to-2-edge model,  $\mathbf{F}^2$ , based on the Singular Value Decomposition (SVD).

### 4.1 A simple proof for universality of second-order tensors

It is enough to consider the  $d_{\text{out}} = 1$  case; the general case is implied by applying the argument for each output feature dimension independently. Let  $\mathbf{G}^2$  be an arbitrary continuous equivariant set-to-2-edge function  $\mathbf{G}^2 : K \subset \mathbb{R}^{n \times d_{\text{in}}} \rightarrow \mathbb{R}^{n \times n}$ . We want to approximate  $\mathbf{G}^2$  with our model  $\mathbf{F}^2$ . First, note that without losing generality we can assume  $\mathbf{G}^2(\mathbf{X})$  has a simple spectrum (i.e., eigenvalues are all different) for all  $\mathbf{X} \in K$ . Indeed, if this is not the case we can always choose  $\lambda > 0$  sufficiently large and consider  $\mathbf{G}^2 + \lambda \text{diag}(1, 2, \dots, n)$ . This diagonal addition does not change

the 2-edge values assigned by  $\mathbf{G}^2$ , and it guarantees a simple spectrum using standard hermitian matrix eigenvalue perturbation theory (see e.g., [32], Section IV:4).

Now let  $\mathbf{G}^2(\mathbf{X}) = \mathbf{U}(\mathbf{X})\mathbf{\Sigma}(\mathbf{X})\mathbf{V}(\mathbf{X})^T$  be the SVD of  $\mathbf{G}^2(\mathbf{X})$ , where  $\mathbf{U} = [\mathbf{u}_1, \dots, \mathbf{u}_n]$  and  $\mathbf{V} = [\mathbf{v}_1, \dots, \mathbf{v}_n]$ . Since  $\mathbf{G}^2(\mathbf{X})$  has a simple spectrum,  $\mathbf{U}, \mathbf{V}, \mathbf{\Sigma}$  are all continuous in  $\mathbf{X}$ ;  $\mathbf{\Sigma}$  is unique, and  $\mathbf{U}, \mathbf{V}$  are unique up to a sign flip of the singular vectors (i.e., columns of  $\mathbf{U}, \mathbf{V}$ ) [24]. Let us first assume that the singular vectors can be chosen uniquely also up to a sign, later we show how we achieve this with some additional mild assumption.

Now, uniqueness of the SVD together with the equivariance of  $\mathbf{G}^2$  imply that  $\mathbf{U}, \mathbf{V}$  are continuous *set-to-set* equivariant and  $\mathbf{\Sigma}$  is a continuous *set invariant* function:

$$\begin{aligned} & (\sigma \cdot \mathbf{U}(\mathbf{X}))\mathbf{\Sigma}(\mathbf{X})(\sigma \cdot \mathbf{V}(\mathbf{X}))^T \\ &= \sigma \cdot \mathbf{G}(\mathbf{X}) = \mathbf{G}(\sigma \cdot \mathbf{X}) \\ &= \mathbf{U}(\sigma \cdot \mathbf{X})\mathbf{\Sigma}(\sigma \cdot \mathbf{X})\mathbf{V}(\sigma \cdot \mathbf{X})^T. \end{aligned} \quad (5)$$

Lastly, since  $\phi$  is set-to-set universal there is a choice of its parameters so that it approximates arbitrarily well the equivariant set-to-set function  $\mathbf{Y} = [\mathbf{U}, \mathbf{V}, \mathbf{1}\mathbf{1}^T\mathbf{\Sigma}]$ . The  $\psi$  component can be chosen by noting that  $\mathbf{G}^2(\mathbf{X})_{i_1, i_2} = \sum_{j=1}^n \sigma_j \mathbf{U}_{i_1, j} \mathbf{V}_{i_2, j} = \mathbf{p}(\beta(\mathbf{Y})_{i_1, i_2, :})$ , where  $\sigma_j$  are the singular values, and  $\mathbf{p} : \mathbb{R}^{6n} \rightarrow \mathbb{R}$  is a cubic polynomial. To conclude pick  $\mathbf{m}$  to approximate  $\mathbf{p}$  sufficiently well so that  $\psi \circ \beta \circ \phi$  approximates  $\mathbf{G}^2$  to the desired accuracy.

To achieve uniqueness of the singular vectors up-to a sign we can add, e.g., the following assumption:  $\mathbf{1}^T \mathbf{u}_i(\mathbf{X}) \neq 0 \neq \mathbf{1}^T \mathbf{v}_i(\mathbf{X})$  for all singular vectors and  $\mathbf{X} \in K$ . Using this assumption we can always pick  $\mathbf{u}_i(\mathbf{X}), \mathbf{v}_i(\mathbf{X})$  in the SVD so that  $\mathbf{1}^T \mathbf{u}_i(\mathbf{X}) > 0, \mathbf{1}^T \mathbf{v}_i(\mathbf{X}) > 0$ , for all  $i \in [n]$ . Lastly, note that equation 5 suggests that also outer-product can be used as a broadcasting layer. We now move to the general proof.

## 4.2 Equivariant set-to- $k$ -edge polynomials

We start with a characterization of the set-to- $k$ -edge equivariant polynomials  $\mathbf{P}^k : \mathbb{R}^{n \times d_{in}} \rightarrow \mathbb{R}^{n^k \times d_{out}}$ . We need some more notation. Given a vector  $\mathbf{x} \in \mathbb{R}^d$ , and a multi-index  $\alpha \in [n]^d$ , we set  $\mathbf{x}^\alpha = \prod_{i=1}^d x_i^{\alpha_i}$ ;  $|\alpha| = \sum_{i=1}^d \alpha_i$ ; and define accordingly  $\mathbf{X}^\alpha = (\mathbf{x}_1^\alpha, \dots, \mathbf{x}_n^\alpha)^T$ . Given two tensors  $\mathbf{A} \in \mathbb{R}^{n^{k_1}}, \mathbf{B} \in \mathbb{R}^{n^{k_2}}$  we use the notation  $\mathbf{A} \otimes \mathbf{B} \in \mathbb{R}^{n^{k_1+k_2}}$  to denote the tensor-product, defined by  $(\mathbf{A} \otimes \mathbf{B})_{i_1, i_2} = \mathbf{A}_{i_1} \mathbf{B}_{i_2}$ , where  $i_1, i_2$  are suitable multi-indices. Lastly, we denote by  $\alpha = (\alpha^1, \dots, \alpha^k)$  a vector of multi-indices  $\alpha^i \in [n]^d$ , and  $\mathbf{X}^\alpha = \mathbf{X}^{\alpha^1} \otimes \dots \otimes \mathbf{X}^{\alpha^k}$ .

**Theorem 3.** *An equivariant set-to- $k$ -edge polynomial  $\mathbf{P}^k : \mathbb{R}^{n \times d_{in}} \rightarrow \mathbb{R}^{n^k \times d_{out}}$  can be written as*

$$\mathbf{P}^k(\mathbf{X}) = \sum_{\alpha} \mathbf{X}^\alpha \otimes \mathbf{q}_\alpha(\mathbf{X}) \quad (6)$$

where  $\alpha = (\alpha^1, \dots, \alpha^k)$ ,  $\alpha^i \in [n]^{d_{in}}$ , and  $\mathbf{q}_\alpha : \mathbb{R}^{n \times d_{in}} \rightarrow d_{out}$  are  $S_n$  invariant polynomials.

As an example, consider the graph case, where  $k = 2$ . Equivariant set-to-2-edge polynomials take the form:

$$\mathbf{P}^k(\mathbf{X}) = \sum_{\alpha^1, \alpha^2} \mathbf{X}^{\alpha^1} \otimes \mathbf{X}^{\alpha^2} \otimes \mathbf{q}_{\alpha^1, \alpha^2}(\mathbf{X}), \quad (7)$$

and coordinate-wise

$$\mathbf{P}_{ijl}^k(\mathbf{X}) = \sum_{\alpha^1, \alpha^2} \mathbf{x}_i^{\alpha^1} \mathbf{x}_j^{\alpha^2} q_{\alpha^1, \alpha^2, l}(\mathbf{X}). \quad (8)$$

The general proof idea and proof itself is given in the supplementary. Figure 3 provides an illustration of these polynomials.

**Approximating  $\mathbf{G}^k$  with a polynomial  $\mathbf{P}^k$ .** We denote for an arbitrary tensor  $\mathbf{A} \in \mathbb{R}^{a \times b \times \dots \times c}$  its infinity norm by  $\|\mathbf{A}\|_\infty = \max_i |\mathbf{A}_i|$ .

**Lemma 1.** *Let  $\mathbf{G}^k : K \subset \mathbb{R}^{n \times d_{in}} \rightarrow \mathbb{R}^{n^k \times d_{out}}$  be a continuous equivariant function over a symmetric domain  $K \subset \mathbb{R}^{n \times d_{out}}$ . For an arbitrary  $\epsilon > 0$ , there exists an equivariant polynomial  $\mathbf{P}^k : \mathbb{R}^{n \times d_{in}} \rightarrow \mathbb{R}^{n^k \times d_{out}}$  so that*

$$\max_{\mathbf{X} \in K} \left\| \mathbf{G}^k(\mathbf{X}) - \mathbf{P}^k(\mathbf{X}) \right\|_\infty < \epsilon.$$

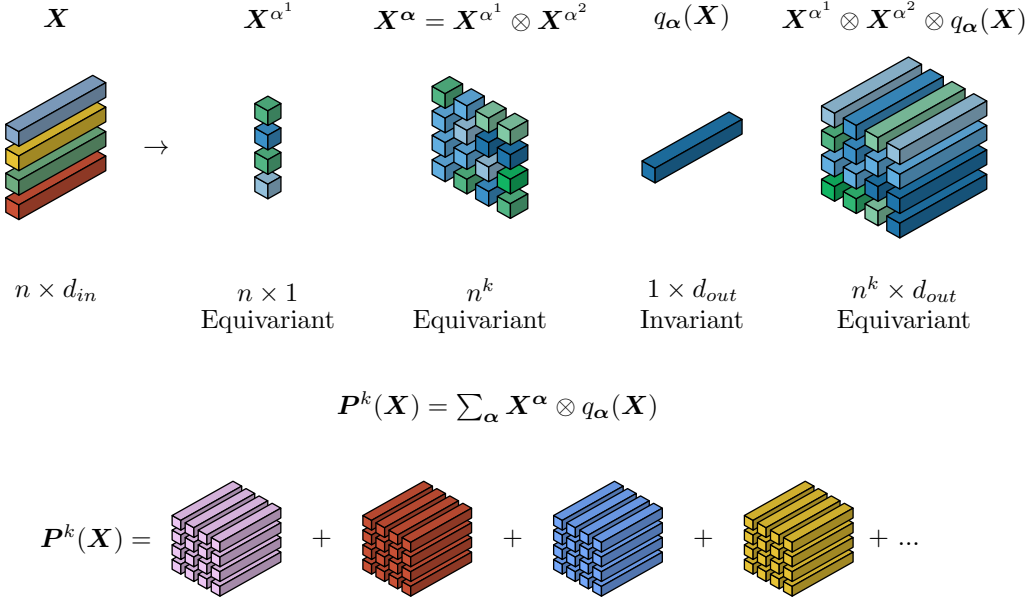


Figure 3: Illustration of the structure of an equivariant set-to- $k$ -edge polynomial  $\mathbf{P}^k$ , for  $k = 2$ .

This is a standard lemma, similar to [42, 21, 30]; we provide a proof in the supplementary.

**Approximating  $\mathbf{P}^k$  with a network  $\mathbf{F}^k$ .** The final component of the proof of Theorem 1 is showing that an equivariant polynomial  $\mathbf{P}^k$  can be approximated over  $K$  using a network of the form in equation 2. The key is to use the characterization of Theorem 3 and write  $\mathbf{P}^k$  in a similar form to our model in equation 2:

$$\mathbf{P}_{i,:}^k(X) = p(\beta(H(X))_{i,:}), \quad (9)$$

where  $H : K \rightarrow \mathbb{R}^{n \times d_1}$  defined by  $H(X)_{i,:} = [x_i, q(X)]$ , where  $q(X) = [q_{\alpha_1}(X), \dots, q_{\alpha_m}(X)]$ , and  $\alpha_1, \alpha_2, \dots, \alpha_m$  are all the multi-indices participating in the sum in equation 6. Note that

$$\beta(H(X))_{i,:} = [x_{i_1}, q(X), x_{i_2}, q(X), \dots, x_{i_k}, q(X)].$$

Therefore,  $p : \mathbb{R}^{d_2} \rightarrow \mathbb{R}^{d_{out}}$  is chosen as the polynomial

$$p : [x_1, \mathbf{y}, x_2, \mathbf{y}, \dots, x_k, \mathbf{y}] \mapsto \sum_{\alpha} x_1^{\alpha_1} \cdots x_k^{\alpha_k} \mathbf{y}_{\alpha},$$

where  $\mathbf{y} = [\mathbf{y}_{\alpha_1}, \dots, \mathbf{y}_{\alpha_m}] \in \mathbb{R}^{m d_{out}}$ , and  $\mathbf{y}_{\alpha_i} \in \mathbb{R}^{d_{out}}$ .

In view of equation 9 all we have left is to choose  $\phi$  and  $\psi$  (i.e.,  $m$ ) to approximate  $H, p$  (resp.) to a desired accuracy. We detail the rest of the proof in the supplementary.

**Universality of the set-to-hypergraph model.** Theorem 2 follows from Theorem 1 by considering a set-to-hypergraph continuous function  $\mathbf{G}$  as a collection  $\mathbf{G}^k$  of set-to- $k$ -edge functions and approximating each one using our model  $\mathbf{F}^k$ . Note that universality still holds if  $\mathbf{F}^1, \dots, \mathbf{F}^K$  all share the  $\phi$  part of the network (assuming sufficient width  $d_1$ ).

Note that a set-to- $k$ -edge model (in equation 2) is not universal when approximating set-to-hypergraph functions:

**Proposition 1.** *The set-to-2-edge model,  $\mathbf{F}^2$ , cannot approximate general set-to-graph functions.*

The proof is in the supplementary; it shows that even the constant function that outputs 1 for 1-edges (nodes), and 0 for 2-edges cannot be approximated by a set-to-2-edge model  $\mathbf{F}^2$ .

## 5 Applications

### 5.1 Model variants and baselines

We tested our model on three learning tasks from two categories: set-to-2-edge and set-to-3-edge.

**Variants of our model.** We consider two variations of our model:

- **S2G:** This is Our basic model. We used the  $\mathbf{F}^2$  and  $\mathbf{F}^3$  (resp.) models for these learning tasks. for  $\mathbf{F}^2$ ,  $\phi$  is implemented using DeepSets [44] with 5 layers and output dimension  $d_1 \in \{5, 80\}$ ;  $\psi$  is implemented with an MLP,  $\mathbf{m}$ , with  $\{2, 3\}$  layers with input dimension  $d_2$  defined by  $d_1$  and  $\beta$ .  $\beta$  is implemented according to equation 3: for  $k = 2$  it uses  $d_2 = 2 * d_1$  output features. For  $\mathbf{F}^3$ , S2G is described in section 5.4.
- **S2G+:** For the  $k = 2$  case we have also tested a more general (but not more expressive) broadcasting  $\beta$  defined using the full equivariant basis  $\mathbb{R}^n \rightarrow \mathbb{R}^{n^2}$  from [20] that contains  $\text{bell}(3) = 5$  basis operations. This broadcasting layer gives  $d_2 = 5 * d_1$ .

**Baselines.** We compare our results to the following baselines:

- **MLP:** A standard multilayer perceptron applied to the flattened set features.
- **SIAM:** A popular similarity learning model (see e.g., [43]) based on Siamese networks. This model has the same structure as in equation 2 where  $\phi$  is a Siamese MLP (a non-universal set-to-set function) that is applied independently to each element in the set. We use the same loss we use with our model (according to the task at hand).
- **SIAM-3:** The same architecture as **SIAM** but with a triplet loss [38] on the learned representations based on  $l_2$  distance, see e.g., [29]. Edge predictions are obtained by thresholding distances of pairs of learned representations.
- **GNN:** A Graph Neural Network [23] applied to the  $k$ -NN ( $k \in \{0, 5, 10\}$ ) induced graph. Edge prediction is done via outer-product [14].
- **AVR:** A non-learnable geometric-based baseline called Adaptive Vertex Reconstruction [36] typically used for the particle physics problem we tackle. More information can be found in the supplementary material.

More architecture, implementation, hyper-parameter details and number of parameters can be found in the supplementary material.

### 5.2 Partitioning for particle physics

The first learning setup we tackle is learning set-to-2-edge functions. Here, each training example is a pair  $(\mathbf{X}, \mathbf{Y})$  where  $\mathbf{X}$  is a set  $\mathbf{X} = (\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n)^T \in \mathbb{R}^{n \times d_{in}}$  and  $\mathbf{Y} \in \{0, 1\}^{n \times n}$  is an adjacency matrix (the diagonal of  $\mathbf{Y}$  is ignored). Our main experiment tackles an important particle partitioning problem.

**Problem statement.** In particle physics experiments, such as the Large Hadron Collider (LHC), beams of incoming particles are collided at high energies. The results of the collision are outgoing particles, whose properties (such as the trajectory) are measured by detectors surrounding the collision point. A critical low-level task for analyzing this data is to associate the particle trajectories to their progenitor, which can be formalized as partitioning sets of particle trajectories into subsets according to their unobserved point of origin in space. This task is referred to as vertex reconstruction in particle physics and is illustrated in Figure 4. We cast this problem as a set-to-2-edge problem by treating the measured particle trajectories as elements in the input set and nodes in the output graph, where the parameters that characterize them serve as the node features. An edge between two nodes indicates that the two particles come from a common progenitor or vertex.

**Data.** We consider three different types (or *flavors*) of particle sets (called *jets*) corresponding to three different fundamental data generating processes labeled bottom-jets, charm-jets, and light-jets (B/C/L). The important distinction between the flavors is the typical number of partitions in each set. Since it is impossible to label real data collected in the detectors at the LHC, algorithms for particle physics are typically designed with high-fidelity simulators, which can provide labeled training data. These algorithms are then applied to and calibrated with real data collected by the LHC experiments. The generated sets are small, ranging from 2 to 14 elements each, with around 0.9M sets divided to

train/val/test using the ratios 0.6/0.2/0.2. Each set element has 10 features ( $d_{in}$ ). More information can be found in the supplementary material.

**Evaluation metrics, loss and post processing.** We consider multiple quantities to quantify the performance of the partitioning: the F1 score, the Rand Index (RI), and the Adjusted Rand Index ( $ARI = (RI - \mathbb{E}[RI]) / (1 - \mathbb{E}[RI])$ ). All models are trained to minimize the F1 score. We make sure the adjacency matrix of the output graph encodes a valid partitioning of nodes to clusters by considering any connected components as a clique.

**Results.** We compare the results of all learning based methods and a typical baseline algorithm used in particle physics (AVR). We also add the results of a trivial baseline that predicts that all nodes have the same progenitor. All models have roughly the same number of parameters. We performed each experiment 11 times with different random initializations, and evaluated the model F1 score, RI and ARI on the test set. The results are shown in Table 2-Jets. For bottom and charm jets, which have secondary vertices, both of our models significantly outperform the baselines by 5%-10% in all performance metrics. In light-jets, without secondary decays, our models yield similar scores. We also performed an extensive ablation study, see Table A1 in the supplementary material. Note that S2G+ has the same expressive power as S2G, and produces equivalent results in practice. Table 1 compares the training times of the different methods.

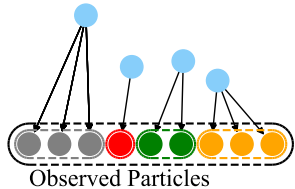


Figure 4: Illustration of a particle physics experiment. The task is to partition the set of observed particles based on their point of origin (in blue).

### 5.3 Learning Delaunay triangulations

In a second set-to-2-edge task we test our model’s ability to learn Delaunay triangulations, namely given a set of planar points we want to predict the Delaunay edges between pairs of points, see e.g., [6] Chapter 9. We generated  $50k$  planar point sets as training data and  $5k$  planar point sets as test data; the point sets,  $\mathbf{X} \in \mathbb{R}^{n \times 2}$ , were uniformly sampled in the unit square, and a ground truth matrix in  $\{0, 1\}^{n \times n}$  was computed per point set using a Delaunay triangulation algorithm. The number of points in a set,  $n$ , is either 50 or varies and is randomly chosen from  $\{20, \dots, 80\}$ . Training was stopped after 100 epochs.

As in the previous experiment, all models have roughly the number of parameters. See more implementation details in the supplementary material. In Table 2-Delaunay we report accuracy of prediction as well as precision recall and F1 score. Evidently, both of our models (S2G and S2G+) outperform the baselines. We also tried the MLP baseline that yielded very low F1 scores (i.e.,  $\leq 0.1$ ). See also Figure 1 in the supplementary material for visualizations of several qualitative examples.

Model	Epochs	training-time (minutes)
S2G	193	62
S2G+	139	47
GNN	91	21
SIAM	77	24
SIAM3	22	322
MLP	132	22

Table 1: Training times of different models. Middle column: number of epochs with early stopping. Right column: total training time in minutes.

### 5.4 Set to 3-edges: learning the convex-hull of a point cloud

In the last experiment, we demonstrate learning of set-to-3-edge function. The learning task we target is finding supporting triangles in the convex hull of a set of points in  $\mathbb{R}^3$ . In this scenario, the input is a point set  $\mathbf{X} \in \mathbb{R}^{n \times 3}$ , and the function we wanted to learn is  $\mathbf{F}^3 : \mathbb{R}^{n \times 3} \rightarrow \mathbb{R}^{n^3}$  where the output is a probability for each triplet of nodes (triangle)  $\{\mathbf{x}_{i_1}, \mathbf{x}_{i_2}, \mathbf{x}_{i_3}\}$  to belong to the triangular mesh that describes the convex hull of  $\mathbf{X}$ . [35] had previously tackled a 2-dimensional version of this problem, but since their network predicts the order of nodes in the 2D convex hull, it is not easily adapted to the 3D settings.

Note that storing 3-rd order tensors in memory is not feasible, hence we concentrate on a *local* version of the problem: Given a point set  $\mathbf{X} \subset \mathbb{R}^3$ , identify the triangles within the  $K$ -Nearest-Neighbors of each point that belong to the convex hull of the entire point set  $\mathbf{X}$ . We used  $K = 10$ . Therefore, for broadcasting ( $\beta$ ) from point data to 3-edge data, instead of holding a 3-rd order tensor in memory we broadcast only the subset of  $K$ -NN neighborhoods. This allows working with high-order information with relatively low memory footprint. Furthermore, since we want to consider 3-edges (triangles)




Model	Jets			Delaunay				Convex Hull (a)			Convex Hull (b)									
	F1	RI	ARI	Acc	Prec	Rec	F1	# points	F1	AUC-ROC	GT	predicted								
B	S2G	0.646±0.003	0.736±0.004	0.491±0.006	$n = 50$				Spherical			$n = 30$								
	S2G+	<b>0.655±0.004</b>	<b>0.747±0.006</b>	<b>0.508±0.007</b>										S2G	0.984	0.927	0.926	0.926	S2G	30
	GNN	0.586±0.003	0.661±0.004	0.381±0.005	S2G+	0.983	<b>0.927</b>	0.925	<b>0.926</b>	GNN5	30	0.693	0.974							
	SIAM	0.606±0.002	0.675±0.005	0.411±0.004	GNN0	0.826	0.384	0.966	0.549	SIAM	30	0.425	0.885							
	SIAM-3	0.597±0.002	0.673±0.005	0.396±0.005	GNN5	0.809	0.363	<b>0.985</b>	0.530	S2G	50	0.686	<b>0.975</b>							
	MLP	0.533±0.000	0.643±0.000	0.315±0.000	GNN10	0.759	0.311	0.978	0.471	GNN5	50	<b>0.688</b>	0.973							
	AVR	0.565	0.612	0.318	SIAM	0.939	0.766	0.653	0.704	SIAM	50	0.424	0.890							
	trivial	0.438	0.303	0.026	SIAM-3	0.911	0.608	0.538	0.570	S2G	20-100	0.535	0.953							
	C	S2G	0.747±0.001	0.727±0.003	0.457±0.004	$n \in \{20, \dots, 80\}$				Gaussian			$n = 50$							
		S2G+	<b>0.751±0.002</b>	<b>0.733±0.003</b>	<b>0.467±0.005</b>										S2G	<b>0.947</b>	<b>0.736</b>	0.934	<b>0.799</b>	S2G
GNN		0.720±0.002	0.689±0.003	0.390±0.005	S2G+	<b>0.947</b>	0.735	0.934	0.798	GNN5	30	0.5826	0.9865							
SIAM		0.729±0.001	0.695±0.002	0.406±0.004	GNN0	0.810	0.387	0.946	0.536	SIAM	30	0.275	0.946							
SIAM-3		0.719±0.001	0.710±0.003	0.421±0.005	GNN5	0.777	0.352	<b>0.975</b>	0.506	S2G	50	<b>0.661</b>	<b>0.997</b>							
MLP		0.686±0.000	0.658±0.000	0.319±0.000	GNN10	0.746	0.322	0.970	0.474	GNN5	50	0.4834	0.9917							
trivial		0.610	0.472	0.078	SIAM	0.919	0.667	0.764	0.687	SIAM	50	0.254	0.974							
AVR		0.695	0.650	0.326	SIAM-3	0.895	0.578	0.622	0.587	S2G	20-100	<b>0.552</b>	<b>0.994</b>							
L		S2G	0.972±0.001	<b>0.970±0.001</b>	<b>0.931±0.003</b>								GNN5	20-100	0.41	0.9866				
		S2G+	0.971±0.002	0.969±0.002	0.929±0.003								SIAM	0.919	0.667	0.764	0.687	SIAM	20-100	0.187
	GNN	0.972±0.001	<b>0.970±0.001</b>	0.929±0.003																
	SIAM	<b>0.973±0.001</b>	<b>0.970±0.001</b>	0.925±0.003																
	SIAM-3	0.895±0.006	0.876±0.008	0.729±0.015																
	MLP	0.960±0.000	0.957±0.000	0.894±0.000																
trivial	0.910	0.867	0.675																	
AVR	0.970	0.965	0.922																	

Table 2: Jets - Performance of partitioning for three types of jets. Delaunay - Results on the Delaunay triangulation task. Convex Hull (a) and (b) - Convex hull learning quantitative and qualitative results.

with no order we used invariant universal set model (DeepSets again) as  $m$ . For  $k = 3$ , S2G is implemented as follows:  $\phi$  is implemented using DeepSets with 3 layers and output dimension  $d_1 = 512$ ;  $\beta$  triplets of points to sets.  $\psi$  is implemented with a DeepSets with 3 layers of 64 features, followed by an MLP,  $m$ , with 3 layers. More details are in the supplementary.

We tested our S2G model on two types of data: Gaussian and spherical. For both types we draw point sets in  $\mathbb{R}^3$  i.i.d. from standard normal distribution,  $\mathcal{N}(0, 1)$ , where for the spherical data we normalize each point to unit length. We generated  $20k$  point set samples as a training set,  $2k$  for validation and another  $2k$  for test set. Point sets are in  $\mathbb{R}^{n \times 3}$ , where  $n = 30$ ,  $n = 50$ , and  $n \in [20, 100]$ . We compare our method, S2G, to the SIAM, GNN and MLP baselines. The F1 scores and AUC-ROC of the predicted convex hull triangles are shown in Table 2-Convex Hull, where our model outperform the baselines in most cases (as in the previous experiment we exclude MLP from the table since it yields very low results). See Figure (b) for several examples of triangles predicted using our trained model compared to the ground truth.

## 5.5 Discussion

In the experiments above we compared our model S2G (and its variant S2G+) with several, broadly used, state-of-the-art architectures. As noted, our models compare favorably with these baselines in all tasks. We attribute this to the increased expressive power (universality) of these models. Our architectures are able to learn relations between set instances while reasoning about the whole set, where Siamese networks rely only on the relation between pairs or triplets of points. The GNN models use a prescribed connectivity that hinder efficient set learning.

## 6 Conclusion

In this paper, we presented a novel neural network family that can model equivariant set-to- $k$ -edge functions and consequently set-to-graph, or more generally set-to-hypergraph functions. The family uses a relatively small number of parameters, compared to other equivariant models, and is shown to be a universal approximator of such continuous equivariant functions. We show the efficacy of these networks on several tasks, including a real-life particle physics problem. There are many directions for future work. One is adapting the model to learn valid clustering of sets (i.e., learn graphs that are built of disjoint cliques). Another direction is incorporating our network architecture in set models, with the goal of improving performance of general set tasks by constructing an intermediate latent graph representation, for example for constructing scene graphs as in [26].

## Broader Impact

Our contribution describes a class of neural network models for functions from sets to graphs and includes theoretical results that the model family is universal. The potential uses of these models is very broad and includes the physical sciences, computer graphics, and social networks. The paper includes experiments that show the positive impact of these models in the context of particle physics, and similar tasks appear repeatedly in the physical sciences. The models could also be used for social networks and areas with more complex ethical and societal consequences. Because the models treat the input as a set and are permutation equivariant, they have the potential to mitigate potential bias in data due to sorting and other pre-processing that could impact methods that treat the input as a sequence. Otherwise, the considerations of bias in the data and impact of failure are no different for our model than the generic considerations of the use of supervised learning and neural networks. Finally we note that the models we describe are already being used in real-world particle physics research.

## Acknowledgments and Disclosure of Funding

HS, NS and YL were supported in part by the European Research Council (ERC Consolidator Grant, "LiftMatch" 771136), the Israel Science Foundation (Grant No. 1830/17) and by a research grant from the Carolito Stiftung (WAIC). JS and EG were supported by the NSF-BSF Grant 2017600 and the ISF Grant 2871/19. KC was supported by the National Science Foundation under the awards ACI-1450310, OAC-1836650, and OAC-1841471 and by the Moore-Sloan data science environment at NYU.

## References

- [1] J. Bromley, I. Guyon, Y. LeCun, E. Säckinger, and R. Shah. Signature verification using a "siamese" time delay neural network. In *Advances in neural information processing systems*, pages 737–744, 1994.
- [2] Z. Chen, S. Villar, L. Chen, and J. Bruna. On the equivalence between graph isomorphism testing and function approximation with gnns. *arXiv preprint arXiv:1905.12560*, 2019.
- [3] T. Cohen and M. Welling. Group equivariant convolutional networks. In *International conference on machine learning*, pages 2990–2999, 2016.
- [4] T. S. Cohen, M. Geiger, J. Köhler, and M. Welling. Spherical cnns. *arXiv preprint arXiv:1801.10130*, 2018.
- [5] T. S. Cohen and M. Welling. Steerable CNNs. (1990):1–14, 2016.
- [6] M. De Berg, M. Van Kreveld, M. Overmars, and O. Schwarzkopf. Computational geometry. In *Computational geometry*, pages 1–17. Springer, 1997.
- [7] S. Dieleman, J. De Fauw, and K. Kavukcuoglu. Exploiting cyclic symmetry in convolutional neural networks. *arXiv preprint arXiv:1602.02660*, 2016.
- [8] C. Esteves, C. Allen-Blanchette, A. Makadia, and K. Daniilidis. 3d object classification and retrieval with spherical cnns. *arXiv preprint arXiv:1711.06721*, 2017.
- [9] J. Gilmer, S. S. Schoenholz, P. F. Riley, O. Vinyals, and G. E. Dahl. Neural message passing for quantum chemistry. In *International Conference on Machine Learning*, pages 1263–1272, 2017.
- [10] Y. Jiang and N. Verma. Meta-learning to cluster. *arXiv preprint arXiv:1910.14134*, 2019.
- [11] N. Keriven and G. Peyré. Universal invariant and equivariant graph neural networks. *CoRR*, abs/1905.04943, 2019.
- [12] T. Kipf, E. Fetaya, K.-C. Wang, M. Welling, and R. Zemel. Neural relational inference for interacting systems. *arXiv preprint arXiv:1802.04687*, 2018.

- [13] T. N. Kipf and M. Welling. Semi-supervised classification with graph convolutional networks. *arXiv preprint arXiv:1609.02907*, 2016.
- [14] T. N. Kipf and M. Welling. Variational graph auto-encoders. *arXiv preprint arXiv:1611.07308*, 2016.
- [15] R. Kondor, H. T. Son, H. Pan, B. Anderson, and S. Trivedi. Covariant compositional networks for learning graphs. *arXiv preprint arXiv:1801.02144*, 2018.
- [16] R. Kondor and S. Trivedi. On the generalization of equivariance and convolution in neural networks to the action of compact groups. *arXiv preprint arXiv:1802.03690*, 2018.
- [17] A. Krizhevsky, I. Sutskever, and G. E. Hinton. Imagenet classification with deep convolutional neural networks. In *Advances in neural information processing systems*, pages 1097–1105, 2012.
- [18] Y. LeCun, L. Bottou, Y. Bengio, P. Haffner, et al. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998.
- [19] H. Maron, H. Ben-Hamu, H. Serviansky, and Y. Lipman. Provably powerful graph networks. *arXiv preprint arXiv:1905.11136*, 2019.
- [20] H. Maron, H. Ben-Hamu, N. Shamir, and Y. Lipman. Invariant and equivariant graph networks. In *International Conference on Learning Representations*, 2019.
- [21] H. Maron, E. Fetaya, N. Segol, and Y. Lipman. On the universality of invariant networks. *arXiv preprint arXiv:1901.09342*, 2019.
- [22] H. Maron, O. Litany, G. Chechik, and E. Fetaya. On learning sets of symmetric elements. *arXiv preprint arXiv:2002.08599*, 2020.
- [23] C. Morris, M. Ritzert, M. Fey, W. L. Hamilton, J. E. Lenssen, G. Rattan, and M. Grohe. Weisfeiler and leman go neural: Higher-order graph neural networks. *arXiv preprint arXiv:1810.02244*, 2018.
- [24] K. A. O’Neil. Critical points of the singular value decomposition. *SIAM journal on matrix analysis and applications*, 27(2):459–473, 2005.
- [25] C. R. Qi, H. Su, K. Mo, and L. J. Guibas. Pointnet: Deep learning on point sets for 3d classification and segmentation. *Proc. Computer Vision and Pattern Recognition (CVPR), IEEE*, 1(2):4, 2017.
- [26] M. Raboh, R. Herzig, J. Berant, G. Chechik, and A. Globerson. Differentiable scene graphs. In *The IEEE Winter Conference on Applications of Computer Vision*, pages 1488–1497, 2020.
- [27] S. Ravanbakhsh, J. Schneider, and B. Póczos. Equivariance through parameter-sharing. *arXiv preprint arXiv:1702.08389*, 2017.
- [28] A. Sannai, Y. Takai, and M. Cordonnier. Universal approximations of permutation invariant/equivariant functions by deep neural networks. *arXiv preprint arXiv:1903.01939*, 2019.
- [29] F. Schroff, D. Kalenichenko, and J. Philbin. Facenet: A unified embedding for face recognition and clustering. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 815–823, 2015.
- [30] N. Segol and Y. Lipman. On universal equivariant set networks. In *International Conference on Learning Representations*, 2020.
- [31] E. Simo-Serra, E. Trulls, L. Ferraz, I. Kokkinos, P. Fua, and F. Moreno-Noguer. Discriminative learning of deep convolutional feature point descriptors. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 118–126, 2015.
- [32] G. W. Stewart. Matrix perturbation theory. 1990.

- [33] P. Veličković, G. Cucurull, A. Casanova, A. Romero, P. Liò, and Y. Bengio. Graph Attention Networks. pages 1–12, 2017.
- [34] O. Vinyals, S. Bengio, and M. Kudlur. Order matters: Sequence to sequence for sets. *arXiv preprint arXiv:1511.06391*, 2015.
- [35] O. Vinyals, M. Fortunato, and N. Jaitly. Pointer networks, 2015.
- [36] W. Waltenberger. RAVE: A detector-independent toolkit to reconstruct vertices. *IEEE Trans. Nucl. Sci.*, 58:434–444, 2011.
- [37] M. Weiler, M. Geiger, M. Welling, W. Boomsma, and T. Cohen. 3D Steerable CNNs: Learning Rotationally Equivariant Features in Volumetric Data. 2018.
- [38] K. Q. Weinberger, J. Blitzer, and L. K. Saul. Distance metric learning for large margin nearest neighbor classification. In *Advances in neural information processing systems*, pages 1473–1480, 2006.
- [39] D. Worrall and G. Brostow. Cubenet: Equivariance to 3d rotation and translation. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 567–584, 2018.
- [40] D. E. Worrall, S. J. Garbin, D. Turmukhambetov, and G. J. Brostow. Harmonic networks: Deep translation and rotation equivariance. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 5028–5037, 2017.
- [41] K. Xu, W. Hu, J. Leskovec, and S. Jegelka. How powerful are graph neural networks? In *International Conference on Learning Representations*, 2019.
- [42] D. Yarotsky. Universal approximations of invariant maps by neural networks. *arXiv preprint arXiv:1804.10306*, 2018.
- [43] S. Zagoruyko and N. Komodakis. Learning to compare image patches via convolutional neural networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4353–4361, 2015.
- [44] M. Zaheer, S. Kottur, S. Ravanbakhsh, B. Poczos, R. R. Salakhutdinov, and A. J. Smola. Deep sets. In *Advances in neural information processing systems*, pages 3391–3401, 2017.