

# DuMLP-Pin: A Dual-MLP-Dot-Product Permutation-Invariant Network for Set Feature Extraction

Jiajun Fei,<sup>1</sup> Ziyu Zhu,<sup>1</sup> Wenlei Liu,<sup>1</sup> Zhidong Deng,<sup>1</sup> Mingyang Li,<sup>2</sup> Huanjun Deng,<sup>2</sup> Shuo Zhang<sup>2</sup>

<sup>1</sup> Institute for Artificial Intelligence at Tsinghua University (THUAI), State Key Laboratory of Intelligent Technology and Systems, Beijing National Research Center for Information Science and Technology (BNRist), Department of Computer Science and Technology, Tsinghua University, Beijing 100084, China

<sup>2</sup> Alibaba Group, Hangzhou 310052, China

feijj20@mails.tsinghua.edu.cn, zhuziyu.edward@gmail.com, {liuw12020, michael}@mail.tsinghua.edu.cn, {mingyangli, huanjun.dhj}@alibaba-inc.com, ygzhangsoya@qq.com

## Abstract

Existing permutation-invariant methods can be divided into two categories according to the aggregation scope, i.e. global aggregation and local one. Although the global aggregation methods, e. g., PointNet and Deep Sets, get involved in simpler structures, their performance is poorer than the local aggregation ones like PointNet++ and Point Transformer. It remains an open problem whether there exists a global aggregation method with a simple structure, competitive performance, and even much fewer parameters. In this paper, we propose a novel global aggregation permutation-invariant network based on dual MLP dot-product, called DuMLP-Pin, which is capable of being employed to extract features for set inputs, including unordered or unstructured pixel, attribute, and point cloud data sets. We strictly prove that any permutation-invariant function implemented by DuMLP-Pin can be decomposed into two or more permutation-equivariant ones in a dot-product way as the cardinality of the given input set is greater than a threshold. We also show that the DuMLP-Pin can be viewed as Deep Sets with strong constraints under certain conditions. The performance of DuMLP-Pin is evaluated on several different tasks with diverse data sets. The experimental results demonstrate that our DuMLP-Pin achieves the best results on the two classification problems for pixel sets and attribute sets. On both the point cloud classification and the part segmentation, the accuracy of DuMLP-Pin is very close to the so-far best-performing local aggregation method with only a 1-2% difference, while the number of required parameters is significantly reduced by more than 85% in classification and 69% in segmentation, respectively. The code is publicly available on <https://github.com/JaronTHU/DuMLP-Pin>.

## Introduction

Convolutional neural networks (CNNs) have superiority in feature extraction on structured data like text, images, or other sequences. For unordered or unstructured set inputs, however, permutation-invariant methods seem to be preferred.

According to the aggregation scope, existing permutation-invariant networks, hereafter referred to as Pin, can be divided into global aggregation and local aggregation methods. Pioneer work uses global aggregation to achieve

permutation-invariance. For instance, PointNet (Qi et al. 2017a) employs global max pooling, while global sum pooling is adopted in Deep Sets (Zaheer et al. 2017). In general, there are simple structures in global aggregation methods. But they are not good at extracting local features. Compared to these, local aggregation methods like PointNet++ (Qi et al. 2017b) and Point Transformer (Zhao et al. 2021) overcome such limitation by defining the local aggregation operation and scope. Accordingly, they yield better performance, though they usually have more complex structures.

In this paper, we propose a dual MLP dot-product permutation-invariant network (DuMLP-Pin) for set feature extraction. The key to DuMLP-Pin is the use of two multilayer perceptrons (MLPs) in order to aggregate information globally. It does gain permutation-invariance through the dot-product of outputs of two MLPs. We strictly prove that the dot-product decomposition always exists for any permutation-invariant function. Compared to other permutation-invariant methods on the same task, DuMLP-Pin is highly parameter-efficient, which means that it only needs a small number of parameters. Furthermore, we also show that the optimization of DuMLP-Pin can be viewed as the constrained optimization of Deep Sets under certain conditions.

The contributions of this paper are summarized as follows:

- We propose a novel global aggregation permutation-invariant network based on dual MLP dot-product decomposition, called DuMLP-Pin.
- We strictly prove that a dot-product decomposition for any permutation-invariant function implemented by DuMLP-Pin always exists if and only if the cardinality of input sets is greater than a pre-specific threshold.
- We explain the high efficiency of DuMLP-Pin and show that under certain conditions DuMLP-Pin can be considered as Deep Sets with strong constraints.

We evaluate DuMLP-Pin on several problems with different types of data elements. Experimental results demonstrate that DuMLP-Pin yields the best performance on both the pixel classification and the attribute set one, compared to other permutation-invariant methods. On the point cloud classification, the performance of the proposed method almost reaches the best-performing local aggregation method

with a difference of only 1.4%, but the number of required parameters is remarkably reduced by more than 85%. In addition, it is also approximately close to the best local aggregation method on the point cloud part segmentation task in performance, where there is only a 1.8% difference between them, but the reduction in the number of required parameters exceeds 69%. Technical appendix is available on <https://arxiv.org/abs/2203.04007>.

## Related Work

### Global Aggregation

Recently, a collection of researches focus on constructing feature extractors for unordered or unstructured data inputs. In pioneer work, global pooling functions are often employed to gain permutation-invariance. Specifically, PointNet (Qi et al. 2017a) utilizes max pooling to directly tackle unordered point cloud data. Feature transformations are also exploited to boost performance in such network architecture. Although PointNet mainly gets involved in point cloud problems, it is easily transferred to other unordered data as well. Meanwhile, Deep Sets (Zaheer et al. 2017; Ravanbakhsh, Schneider, and Poczos 2017) makes use of sum pooling to aggregate information. It also points out that conditional mapping could be adopted as additional meta-information is available. The core of Deep Sets is to add up all representations and then apply nonlinear transformations. BRUNO (Korshunova et al. 2018) is a recurrent model for Bayesian inference on set inputs, where joint distribution over observations is permutation-invariant. Besides the above-mentioned simple pooling functions, attention-based architecture is also popular when building permutation-invariant or permutation-equivariant functions. Ilse, Tomczak, and Welling proposes an attention-based aggregation operator to learn the Bernoulli distribution related to multiple instance learning problems (Maron and Lozano-Pérez 1998). Set Transformer (Lee et al. 2019) is another Transformer-based neural network that is designed to process sets of data. In such a model, several attention-related set operations are presented under different problem settings. There are many research works associated with permutation-invariance and permutation-equivariance in Transformer (Zhao, Jia, and Koltun 2020; Khan et al. 2021), where some of them leverage permutation-invariance, but some do not. Unlike previous methods directly approximating global features, RepSet (Skianis et al. 2020) generates them by solving the maximum matching problem between the elements of two sets.

Global aggregation methods have simple structures, whose permutation-invariance is gained by either global pooling functions or attention mechanisms. Such two types of methods, however, have different limitations. Wagstaff et al. analyzes the limitation of pooling architectures and provides a solution by imposing a constraint on the dimensionality of the latent space. For Transformer-based global aggregations, there is usually a relatively heavy computational complexity. Some researchers argue that full attention in Transformer is not necessary (Sukhbaatar et al. 2019; Kitaev, Kaiser, and Levskaya 2019; Liu et al. 2021), and accordingly, different ways are delivered to do computation

more efficiently.

### Local Aggregation

Global aggregation methods are always weak in capturing local information, which generally results in relatively poor performance, especially in the point cloud settings. To overcome this drawback, local aggregation methods are proposed. PointNet++ (Qi et al. 2017b) captures local structures by recursively applying PointNet on the partitioning of point cloud inputs. In the PointNet++ model, two density adaptive layers, i.e. multi-scale grouping and multi-resolution grouping, are proposed to group local regions and combine features from different scales. PointASNL (Yan et al. 2020) focuses on tackling point clouds with noise and presents adaptive sampling and local-nonlocal modules to seize local and long-range dependencies. Point Cloud Transformer (Guo et al. 2021) is also based on Transformer architectures. Unlike Set Transformer, however, it uses the offset-attention module with neighbor embedding for augmented local feature aggregation. Point Transformer (Zhao et al. 2021) establishes attention mechanisms based on (Zhao, Jia, and Koltun 2020). It also employs relative positional encoding to aggregate local information. In summary, local aggregation methods generally have more complex structures since they have to carefully define aggregation operations and scopes.

### DuMLP-Pin

In this section, we first formulate the permutation-invariant learning problem and then propose a novel global aggregation permutation-invariant network based on dual MLP dot-product decomposition (DuMLP-Pin).

### Problem Formulation

Suppose that  $P$  denotes the permutation matrix and  $\mathcal{P}_N$  the set of all permutation matrices with size  $N \times N$ , where  $N$  is the number of elements in a given input set. Let  $p$  and  $r$  indicate the dimension of the input and output vectors, respectively. According to the definition of permutation-invariance and permutation-equivariance (Zaheer et al. 2017), we can rewrite them in a matrix form.

**Definition 1 (Permutation-invariant)** *The function  $f : \mathbb{R}^{N \times p} \rightarrow \mathcal{Y}$  is permutation-invariant if*

$$f(PX) = f(X) \text{ for } \forall X \in \mathbb{R}^{N \times p}, \forall P \in \mathcal{P}_N. \quad (1)$$

**Definition 2 (Permutation-equivariant)** *The function  $g : \mathbb{R}^{N \times p} \rightarrow \mathbb{R}^{N \times r}$  is permutation-equivariant if*

$$g(PX) = Pg(X) \text{ for } \forall X \in \mathbb{R}^{N \times p}, \forall P \in \mathcal{P}_N. \quad (2)$$

In this paper, we pay attention to the set classification problem with supervised learning. Accordingly, the permutation-invariant learning problem can be described as follows:

**Definition 3 (Permutation-invariant Learning)** *Given  $M$  training samples in the form of  $\{(X_i, y_i)\}_{i=1}^M$ , where  $X_i \in \mathbb{R}^{N \times p}$  represents a set input and  $y_i \in \mathcal{Y}$  a label, find the function  $f : \mathbb{R}^{N \times p} \rightarrow \mathcal{Y}$  such that the corresponding  $f$ -related performance metric evaluated on test sets is maximized.*

## Minimum Dimension Decomposition

We can construct a bijection  $u : \mathbb{R}^{m \times n} \rightarrow \mathbb{R}^{mn}$  through simply reshaping, which means that  $u(x)_{in+j} = x_{i,j}$ . Clearly, if a function can represent any matrices of size  $m \times n$ , it can then express any vectors of length  $mn$  by adding a flattening module. In the following discussion, we will focus on the representation of feature matrices instead of feature vectors. First of all, we introduce Lemma 1 for minimum dimension decomposition (MDD).

**Lemma 1 (MDD)**  $\forall A \in \mathbb{R}^{m \times n}, m \leq n, \forall l \geq m, \forall B \in \{X \in \mathbb{R}^{m \times l} \mid \text{rank}(X) = m\}$ . Let  $S = \{C \in \mathbb{R}^{l \times n} \mid BC = A\}$ , then

$$S \neq \emptyset \text{ and } S = \left\{ C_p + X_h \Lambda \mid \Lambda \in \mathbb{R}^{(l-m) \times n} \right\}, \quad (3)$$

where  $C_p$  stands for a particular solution of  $BC = A$ , and the columns of  $X_h$  are basis of  $\text{Ker}(B)$ .

The proof is supplemented in the technical appendix. Obviously there exist infinite tuples  $(B, C)$  for any  $A$  through MDD. Therefore, we have more choices when we represent  $A$  through the decomposition of  $B$  and  $C$ .

## Permutation-Invariance through Dot-Product Decomposition

By substituting matrices into matrix functions in Lemma 1, we may always decompose one matrix function by two matrix functions. In Theorem 1, we show that any permutation-invariant function can be decomposed into two permutation-equivariant functions through a dot-product decomposition.

**Theorem 1** For any permutation-invariant function  $f : \mathbb{R}^{N \times p} \rightarrow \mathbb{R}^{s \times t}$ , there exist two permutation-equivariant functions  $g^{(1)} : \mathbb{R}^{N \times p} \rightarrow \mathbb{R}^{N \times s}, g^{(2)} : \mathbb{R}^{N \times p} \rightarrow \mathbb{R}^{N \times t}$ , such that

$$\forall X \in \mathbb{R}^{N \times p}, f(X) = \left( g^{(1)}(X) \right)^T g^{(2)}(X) \quad (4)$$

if and only if  $N \geq \min\{s, t\}$ .

The proof is also seen in the technical appendix for details. Unlike the theorems given in previous work (Qi et al. 2017a; Zaheer et al. 2017), Theorem 1 discuss the relationship between the shape of the feature matrix  $\min\{s, t\}$  and the cardinality of the input set  $N$ . An intuition is that when  $N < \min\{s, t\}$ , the term  $(g^{(1)}(X))^T g^{(2)}(X)$  cannot be full-rank according to linear algebra. Therefore, it cannot represent full-rank feature matrices, which extremely limits the expression capability. In practice, it is much easier to construct a permutation-equivariant function since any function that processes each element independently is permutation-equivariant. It is also the main reason that we refer to two permutation-equivariant functions to represent one permutation-invariant function. When generating feature vectors of length  $C$ , the outputs of both MLPs can be only as small as  $\sqrt{C}$  (when  $N \geq \sqrt{C}$ ), which leads to a great reduction of parameters.

## Aggregation Block and Broadcast Block

DuMLP-Pin is mainly constructed by two kinds of blocks: aggregation block and broadcast block. The aggregation block, which performs global feature aggregation, consists of two vanilla MLPs. It achieves permutation-invariance through the dot-product of the outputs of these MLPs. However, sometimes we need not only set features that represent all unordered or unstructured elements, but also element features for element label prediction. Thus we need to broadcast the set feature to each element. In the following, we design the broadcast block to address such a problem. Formally, let element feature  $x \in \mathbb{R}^{d_x}$ , set feature  $y \in \mathbb{R}^{d_y}$ . Then we use Equation 5 to generate combined feature  $z \in \mathbb{R}^{d_z}$ .

$$z = (W_{xy} \cdot x) \cdot y + W_x \cdot x + W_y \cdot y + b, \quad (5)$$

where  $W_{xy} \in \mathbb{R}^{d_z \times d_y \times d_x}, W_x \in \mathbb{R}^{d_z \times d_x}, W_y \in \mathbb{R}^{d_z \times d_y}$ , and  $b \in \mathbb{R}^{d_z}$ . Although  $W_{xy}$  involves direct interaction between  $x$  and  $y$ , it increases computational burden tremendously. We observe that involving such term only leads to insignificant performance improvement in experiment. Consequently, we balance efficiency and effectiveness by setting  $W_{xy} = 0$ . The structures of aggregation block and broadcast block are shown in Figure 1 (a) and (b), respectively.

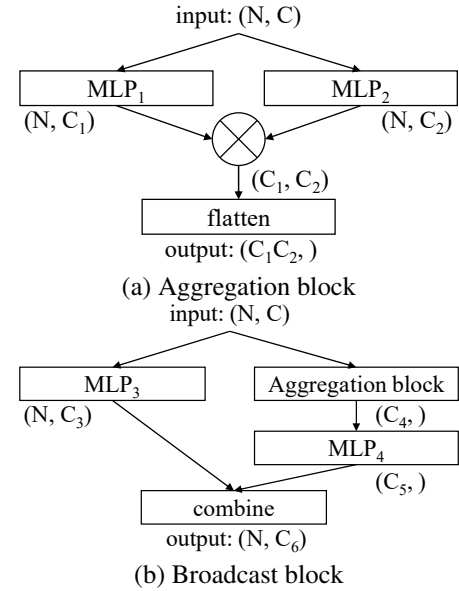


Figure 1: Structure design for aggregation block and broadcast block.

## Nonlinear Activation Function

Nonlinear activation functions are essential for the success of deep learning. Besides the common-used ReLU (Nair and Hinton 2010), there are many works focusing on designing novel nonlinear activation functions, including ELU (Clevert, Unterthiner, and Hochreiter 2016), SELU (Klambauer et al. 2017), Mish (Misra 2020), etc. We do not discuss general preferences on the choice of nonlinear activation functions. Instead, we only focus on the last layers of two MLPs in the aggregation block.

Without loss of generality, suppose that both MLPs comprise only one linear layer without any bias. Let  $W_1$  and  $W_2$  be weights.  $\alpha_1$  and  $\alpha_2$  are assumed to be nonlinear activation functions used by MLP<sub>1</sub> and MLP<sub>2</sub>, respectively. Since we argue both MLPs must be permutation-equivariant, here we only discuss element-wise activation functions like ReLU and set-wise ones like softmax. Then for input  $X \in \mathbb{R}^{N \times C}$ , the output is computed as:

$$f(x) = \alpha'_1(W_1^T X^T) \alpha_2(X W_2), \quad (6)$$

where  $\alpha'_1$  satisfies that  $\alpha'_1(W_1^T X^T) = \alpha_1^T(X W_1)$ . We find that the model always shows better performance when using at least one nonlinear activation function. Why nonlinear activation functions are essential? Assume that both MLPs do not use nonlinear activation function. In this case, Equation 6 turns into Equation 7, i.e.,

$$f(x) = W_1^T X^T X W_2 = W_1^T (X^T X) W_2. \quad (7)$$

It is easy to see that  $f(X)$  is dominated by  $X^T X \in \mathbb{R}^{C \times C}$ . This is inappropriate and even destructive especially when  $C \ll N$ . For example, the feature dimension of the input element is relatively small under some point cloud settings, which is much smaller than the number of points. If no nonlinear activation function is used, then unfeasible constraints that for  $X_1, X_2 \in \mathbb{R}^{N \times C}$ ,  $X_1^T X_1 = X_2^T X_2 \Rightarrow f(X_1) = f(X_2)$  will be imposed. However, using even one nonlinear activation function can overcome this problem by cutting off the direct combination between  $X^T$  and  $X$ .

Besides, we observe that set-wise softmax generally outperforms element-wise activation functions. We argue that set-wise softmax involves direct interaction among elements in inputs, thus leading to better results than element-wise ones. We will discuss this more in the ablation study.

## Further Discussion

In this part, we methodologically make comparisons of our DuMLP-Pin with other global aggregation methods including Deep Sets, PointNet, and Set Transformer.

**Relation to both Deep Sets (Zaheer et al. 2017) and PointNet (Qi et al. 2017a).** Both Deep Sets and PointNet build permutation-invariant functions in a similar manner:

$$\text{Deep Sets: } f(X) = \gamma \left( \sum_{x_i} \{h(x_i)\} \right),$$

$$\text{PointNet: } f(X) = \gamma \left( \max_{x_i} \{h(x_i)\} \right).$$

The main difference lies in the pooling function, which is also the core of permutation-invariance. To illustrate how DuMLP-Pin works, let us suppose that each row is processed independently. It implies that any interaction between elements does not get involved. Thus we make the derivation below:

$$\begin{aligned} f(X) &= \left( g^{(1)}(X) \right)^T g^{(2)}(X) \\ &= \sum_{x_i} \left( g^{(1)}(x_i) \right)^T \left( g^{(2)}(x_i) \right) \\ &= \sum_{x_i} h'(x_i). \end{aligned}$$

Therefore, DuMLP-Pin also takes advantage of the sum pooling over element features. However, unlike  $h(x_i)$  in Deep Sets,  $h'(x_i) \in \mathbb{R}^{C_1 \times C_2}$  is subject to the following strong constraints:

$$\text{rank} [h'(x_i)] \leq 1. \quad (8)$$

From this perspective, DuMLP-Pin can be regarded as Deep Sets with strong constraints on element features. We argue that the direct optimization of  $h$  is generally tougher and more unstable than optimizing  $g^{(1)}$  and  $g^{(2)}$  concurrently. This is because we shrink the search space by imposing constraints without loss of theoretical expression capability. If we use set-wise softmax, then  $g(x_i)$  can be modified as  $g(x_i)p(X)$  where  $p(X)$  is a global pooling. But it does not change the low-rank constraints.

**Relation to Set Transformer (Lee et al. 2019).** Although Pooling by Multihead Attention Block (PMA) in Set Transformer is similar to our aggregation block, they are different in many aspects as shown in Equation 9. Set Transformer is based on Transformer architecture, which implies that it mainly employs linear mapping for feature transformation, while our DuMLP-Pin uses MLPs for the same sake. In addition, the aggregation block does not leverage the self-attention mechanism, layer normalization, and multi-head architecture. The output of Set Transformer has different meanings on different dimensions (i.e. sample and feature). But our DuMLP-Pin takes both dimensions as features, making it natural to flatten for further use.

$$f_{\text{PMA}}(X) = \text{Attention}(S, X, X),$$

$$f_{\text{DuMLP-Pin}}(X) = \text{Flatten} \left[ \text{MLP}_1(X)^T \text{MLP}_2(X) \right]. \quad (9)$$

## Possibility of High-order Dot-Product Decomposition.

DuMLP-Pin can be viewed as Deep Sets with strong constraints on element features, which is derived by raising feature dimensions from 1D vector to 2D matrix. Is it possible that we construct 3D, 4D, or even  $n$ D feature tensors to further shrink search space and boost performance? We dive into such a problem and formulate the result as Theorem 2.

**Theorem 2** *For any permutation-invariant function  $f : \mathbb{R}^{N \times p} \rightarrow \mathbb{R}^{c_1 \times \dots \times c_n}$ , there exist  $n$  permutation-equivariant functions  $g^{(1)} : \mathbb{R}^{N \times p} \rightarrow \mathbb{R}^{N \times c_1}, \dots, g^{(n)} : \mathbb{R}^{N \times p} \rightarrow \mathbb{R}^{N \times c_n}$ , such that ( $1 \leq a_j \leq c_j$ )*

$$\forall X \in \mathbb{R}^{N \times p}, f(X)_{a_1, \dots, a_n} = \sum_{i=1}^N \prod_{j=1}^n \left( g^{(j)}(X) \right)_{i, a_j}$$

$$\text{if } N \geq \frac{\prod_{j=1}^n c_j}{\max_j c_j} \text{ and only if } N \geq \frac{\prod_{j=1}^n c_j}{\sum_{j=1}^n c_j}$$

$$\text{and } N \geq \max_{C \subset \{c_1, \dots, c_n\}} \min \left\{ \prod_{c \in C} c, \prod_{c \in \bar{C}} c \right\}.$$

The proof is also supplemented in the technical appendix. We compare the performance of different feature orders in the ablation study.

## Experiments

On several problems with different types of data elements, we evaluate the performance of our DuMLP-Pin. In the experiments, we employ the MNIST for pixel set classification, the CelebA for attribute set anomaly detection, the ModelNet40 for point cloud classification, and the ShapeNetPart for point cloud part segmentation, respectively. We independently run all experiments 5 times and report the evaluation metric on test sets. More details about the architectures and hyper-parameters are available in the technical appendix.

### Pixel Set Classification

The MNIST (LeCun et al. 1998) is a well-known handwritten digit classification dataset with 60,000 images in the training set and 10,000 images in the test set. We transform images into pixel sets by appending relative coordinates to intensities and shuffling the order randomly. Accordingly, our DuMLP-Pin is compared with PointNet and PointNet++, respectively, since they also consider images as pixel sets. We use the error rate as the evaluation metric. Both results quoted are taken from the cited papers.

Method	#Params(M)	Error Rate(%)
PointNet (2017a)	3.47	0.78
PointNet++ (2017b)	1.47	0.51
DuMLP-Pin-S	0.28	0.80 ± 0.053
DuMLP-Pin-L	1.13	<b>0.48</b> ± 0.034

Table 1: The MNIST pixel set classification results.

We evaluate two models of different sizes: DuMLP-Pin-S with only one aggregation block and DuMLP-Pin-L with one aggregation block and two broadcast blocks. From Table 1, it is readily observed that DuMLP-Pin-S performs slightly worse than PointNet while DuMLP-Pin-L is slightly better than PointNet++. However, unlike PointNet++ as a local aggregation method, DuMLP-Pin has no access to local structures, making the results more competitive. It is possible to construct a DuMLP-Pin model with better hyper-parameters, which may further improve the performance. The experiment shows that DuMLP-Pin has the potential to process simple images. We use t-SNE (van der Maaten and Hinton 2008) to visualize feature vectors preceding classifiers, as shown in Figure 2.

### Attribute Set Anomaly Detection

CelebA (Liu et al. 2015) is a large-scale face attributes dataset including 202,599 face images with 40 attribute annotations per image. We randomly sample 9 images with the same 2 attributes and 1 image without those. Each image in CelebA is forced to appear at most once in the sampled dataset. We formulate this problem as a set feature extraction problem because the anomaly is in the minority and most normal members do not deviate much from the set feature. All methods are the same except the set feature extractor part.

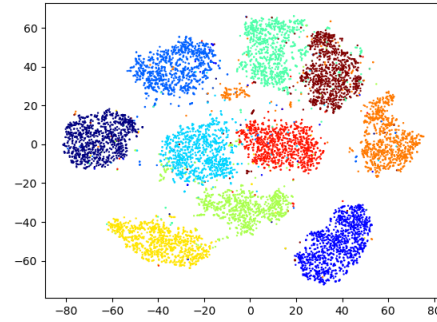


Figure 2: The t-SNE embeddings of DuMLP-Pin-S on the MNIST test set. Feature vectors are generated by one aggregation block with only 17.9K parameters.

Method	#Params(K)	OA(%)
Max-pooling	65.8	60.4 ± 0.44
Mean-pooling	65.8	65.0 ± 0.16
PointNet (2017a)	131.6	62.3 ± 0.22
Deep Sets (2017)	132.6	65.0 ± 0.35
Set Transformer (2019)	263.4	66.4 ± 0.45
DuMLP-Pin	83.3	<b>66.8</b> ± 0.43

Table 2: The CelebA attribute set anomaly detection results.

We evaluate all methods under the end-to-end setting without access to the attributes. We use overall accuracy (OA) as the evaluation metric. The structures of other methods are selected from a series of candidates through grid search. From Table 2, DuMLP-Pin and Set Transformer outperform other methods by a significant margin, but DuMLP-Pin using fewer parameters. Figure 3 shows some samples from the set anomaly detection dataset.



Figure 3: Samples from set anomaly detection datasets, the blue frame denotes anomaly and the red frame denotes wrong-classified image. We visualize the probability by adding a blue bar under each image.

### Point Cloud Classification

The ModelNet40 (Wu et al. 2015) is a 40-class 3D shape dataset with 9,843 CAD models in the training set and 2,468 CAD models in the test set, respectively. We utilize 1,024

points with normal vectors from the sampled point cloud of Qi et al. for fair comparisons. We take overall accuracy as the evaluation metric. All results quoted are taken from the cited papers.

Method	#Params(M)	OA(%)
PointNet (2017a)	3.48	89.2
Deep Sets (2017)	1.34	90 ± 0.3
Set Transformer (2019)	1.14	90.4 ± 1.73
PointNet++ (2017b)	1.74	91.9
PCT (2021)	2.88	93.2
Point Transformer (2021)	~9.6*	<b>93.7</b>
DuMLP-Pin	<b>0.29</b>	92.3 ± 0.12

Table 3: The ModelNet40 point cloud classification results. Note that \* indicates that the number of parameters is estimated from unofficial code implementation.

We also evaluate two models of different sizes, but they have similar performance, meaning that the model is prior to overfitting when adding more parameters. Therefore, we only focus on the small one. The performance of DuMLP-Pin is better than other methods for sets (Deep Sets and Set Transformer) and pioneer methods for point clouds (PointNet and PointNet++), but still 1%~2% worse than state-of-the-art local aggregation methods for point clouds. As a global aggregation method, DuMLP-Pin cannot capture more local information. But the result is competitive since the structure of DuMLP-Pin is simpler and the number of required parameters is remarkably decreased by more than 85%. We adopt the same model retrieval method as that introduced in (Qi et al. 2017a) in Figure 4.

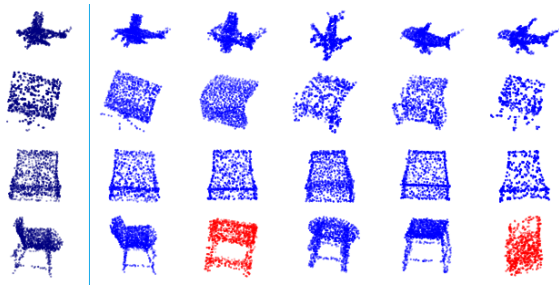


Figure 4: Model retrieval results from the ModelNet40 test set. The left column is the query point cloud, and the others are top-5 retrieved point clouds. Red means incorrect classification.

### Point Cloud Part Segmentation

The ShapeNetPart (Yi et al. 2016) is a 3D object part segmentation dataset with 14,006 objects in the training set and 2,874 ones in the test set. It consists of 16 shape categories and 50 part ones, respectively. We utilize 2,048 points like the sampled point clouds of Qi et al. for fair comparisons. The mIoU is employed as the evaluation metric and all their results are quoted from the cited papers.

Method	#Params(M)	mIoU(%)
PointNet (2017a)	7.83	83.7
PointNet++ (2017b)	1.30	85.1
PCT (2021)	3.14	86.4
Point Transformer (2021)	~19.2*	<b>86.6</b>
DuMLP-Pin-S	<b>0.70</b>	83.5 ± 0.06
DuMLP-Pin-L	0.98	84.8 ± 0.04

Table 4: The ShapeNetPart point cloud part segmentation results. Note that \* implies that the data is roughly calculated based on unofficial code implementation.

Like that in previous experiments, we evaluate two models of different sizes: DuMLP-Pin-S with two broadcast blocks and DuMLP-Pin-L with four broadcast blocks. From Table 4, it is observed that DuMLP-Pin-S somewhat underperforms PointNet and DuMLP-Pin-L is slightly inferior to PointNet++. But DuMLP-Pin has its own advantages. For instance, it does not use any normal vectors as input. Without the use of local features, it is reasonable that the performance of DuMLP-Pin is only 2% lower than that of attention-based local aggregation methods for point cloud problems. We visualize some of the segmentation examples obtained by DuMLP-Pin in Figure 5.

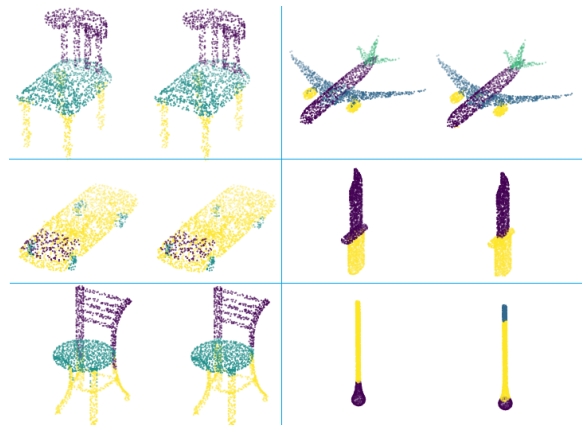


Figure 5: Segmentation results on the ShapeNetPart test set. For each group of point clouds, on the left is the ground truth, while on the right is our prediction.

### Computational Complexity Analysis

We compare the computational complexity of different methods in the point cloud classification task, as listed in Table 5. It is readily observed that DuMLP-Pin has the lowest number of parameters and MACs per sample among all the methods. There are two reasons behind such minimum time-consuming advantage. First, the proposed method has smaller hidden layer dimensions. Second, it also has linear complexity with respect to the number of points, while the computational complexity of transformer-based methods is generally quadratic.

Method	#Params(M)	MACs/sample(M)
PointNet	3.48	440
Deep Sets	1.34	135
Set Transformer	1.14	1,284
PointNet++	1.74	3,932
PCT	2.88	2,286
Point Transformer	$\sim 9.6^*$	$\sim 18,314^*$
DuMLP-Pin	<b>0.29</b>	<b>19</b>

Table 5: The comparison of time efficiency among all the methods in the point cloud classification. \* indicates that the value is estimated from unofficial code implementation.

## Ablation Study

The analysis of the influence of several factors including nonlinear activation function, factorization, and decomposition order are further conducted. All experiments are done with ModelNet40, where all MLPs adopted in DuMLP-Pin contain three hidden layers with sizes of 32, 128, and  $k$  which may vary in different experiments, respectively.

**Nonlinear Activation Function.** All models are the same except the activation functions of two MLPs. For set-wise activation functions, we adopt squashing introduced in (Sabour, Frosst, and Hinton 2017) besides softmax. Figure 6 visualizes different activation functions used in our work. It can be seen that softmax dominates the table with all metrics around 1% greater than their non-softmax counterparts. Squashing seems to lose so much information that it has poor performance. Generally, we prefer using two softmax or one softmax and the other no activation.

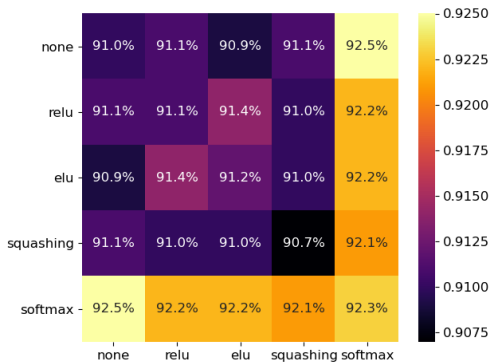


Figure 6: Influence of nonlinear activation function.

**Factorization.** By controlling the dimension of feature vectors as 1,024, we experiment on several different factorizations. The influence of factorization does not seem critical because the OA of the best model is only 0.7% larger than the worst model in Table 6. However, since the square factorization ( $32 \times 32$ ) needs the fewest parameters, it is highly recommended to use square factorization first and try other factorizations as well if possible. To control the variable, we use both activation functions as softmax, which makes the

result slightly different from that with one single softmax in Figure 6.

Factorization	#Params(K)	OA(%)
$1 \times 1,024$	143.8	91.8
$2 \times 512$	76.9	92.1
$4 \times 256$	43.6	91.6
$8 \times 128$	27.4	91.7
$16 \times 64$	20.0	91.9
$32 \times 32$	<b>17.9</b>	<b>92.3</b>

Table 6: Influence of factorization.

**Decomposition Order.** By setting  $N = 1,024$  and controlling the dimension of feature vectors as 1,024, the decomposition always exists according to Theorem 2. Because multiple softmax makes the value too small, we only use one softmax in all experiments except the first one (since it all degenerates to one). From Table 7, the performance of order 2 is better than all the other methods. The possible reason is that the use of high-order decomposition is equivalent to adding stronger constraints. In this case, it is hard to use any vanilla MLPs to parameterize each function in the decomposition.

Order	Factorization	#Params(K)	OA(%)
1	1,024	138.9	90.3
2	$32^2$	<b>17.9</b>	<b>92.5</b>
3	$16 \times 8^2$	18.5	92.0
4	$8^2 \times 4^2$	22.2	91.6
5	$4^5$	26.5	91.2
10	$2^{10}$	50.3	91.3

Table 7: Influence of decomposition order.

## Conclusion

In this paper, we propose a novel global aggregation permutation-invariant network based on dual MLP dot-product decomposition, which is called DuMLP-Pin. We strictly prove that two or more permutation-equivariant functions can be leveraged to set up one permutation-invariant function. The high efficiency of DuMLP-Pin is also analyzed by converting the optimization problem of DuMLP-Pin to that of Deep Sets with strong constraints. Finally, we evaluate the performance of our DuMLP-Pin on different tasks. The experimental results illustrate that DuMLP-Pin achieves the best results on the pixel set classification and the attribute set anomaly detection. Additionally, the performance of DuMLP-Pin is only slightly lower than existing SOTA methods, while the number of required parameters is significantly reduced. In the future, we hope the DuMLP-Pin could extensively be exploited as a competitive global aggregation baseline for a wide range of unordered or unstructured data problems in terms of its distinguishing performance.

## Acknowledgments

This work was supported by Alibaba Group through Alibaba Innovative Research Program. It was also supported in part by the National Science Foundation of China (NSFC) under Grant No. 62176134 and by a grant from the Institute Guo Qiang (2019GQG0002), Tsinghua University.

## References

- Clevert, D.; Unterthiner, T.; and Hochreiter, S. 2016. Fast and Accurate Deep Network Learning by Exponential Linear Units (ELUs). In Bengio, Y.; and LeCun, Y., eds., *Proceedings of the 4th International Conference on Learning Representations*.
- Guo, M.-H.; Cai, J.-X.; Liu, Z.-N.; Mu, T.-J.; Martin, R. R.; and Hu, S.-M. 2021. PCT: Point cloud transformer. *Computational Visual Media*, 7(2): 187–199.
- Ilse, M.; Tomczak, J.; and Welling, M. 2018. Attention-based Deep Multiple Instance Learning. In *Proceedings of the 35th International Conference on Machine Learning*, 2127–2136. PMLR.
- Khan, S.; Naseer, M.; Hayat, M.; Zamir, S. W.; Khan, F. S.; and Shah, M. 2021. Transformers in Vision: A Survey. arXiv:2101.01169.
- Kitaev, N.; Kaiser, L.; and Levskaya, A. 2019. Reformer: The Efficient Transformer. In *Proceedings of the 7th International Conference on Learning Representations*.
- Klambauer, G.; Unterthiner, T.; Mayr, A.; and Hochreiter, S. 2017. Self-Normalizing Neural Networks. In Guyon, I.; Luxburg, U. V.; Bengio, S.; Wallach, H.; Fergus, R.; Vishwanathan, S.; and Garnett, R., eds., *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc.
- Korshunova, I.; Degraeve, J.; Huszar, F.; Gal, Y.; Gretton, A.; and Dambre, J. 2018. BRUNO: A Deep Recurrent Model for Exchangeable Data. In Bengio, S.; Wallach, H.; Larochelle, H.; Grauman, K.; Cesa-Bianchi, N.; and Garnett, R., eds., *Advances in Neural Information Processing Systems*, volume 31. Curran Associates, Inc.
- LeCun, Y.; Bottou, L.; Bengio, Y.; and Haffner, P. 1998. Gradient-Based Learning Applied to Document Recognition. *Proceedings of the IEEE*, 86(11): 2278–2324.
- Lee, J.; Lee, Y.; Kim, J.; Kosiorek, A.; Choi, S.; and Teh, Y. W. 2019. Set Transformer: A Framework for Attention-based Permutation-Invariant Neural Networks. In *Proceedings of the 36th International Conference on Machine Learning*, 3744–3753. PMLR.
- Liu, Z.; Lin, Y.; Cao, Y.; Hu, H.; Wei, Y.; Zhang, Z.; Lin, S.; and Guo, B. 2021. Swin Transformer: Hierarchical Vision Transformer Using Shifted Windows. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 10012–10022.
- Liu, Z.; Luo, P.; Wang, X.; and Tang, X. 2015. Deep Learning Face Attributes in the Wild. In *Proceedings of the IEEE International Conference on Computer Vision*, 3730–3738.
- Maron, O.; and Lozano-Pérez, T. 1998. A Framework for Multiple-Instance Learning. In Jordan, M.; Kearns, M.; and Solla, S., eds., *Advances in Neural Information Processing Systems*, volume 10. MIT Press.
- Misra, D. 2020. Mish: A Self Regularized Non-Monotonic Neural Activation Function. In *Proceedings of the 31st British Machine Vision Virtual Conference*.
- Nair, V.; and Hinton, G. E. 2010. Rectified Linear Units Improve Restricted Boltzmann Machines. In *Proceedings of the 27th International Conference on Machine Learning*.
- Qi, C. R.; Su, H.; Mo, K.; and Guibas, L. J. 2017a. PointNet: Deep Learning on Point Sets for 3D Classification and Segmentation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*.
- Qi, C. R.; Yi, L.; Su, H.; and Guibas, L. J. 2017b. PointNet++: Deep Hierarchical Feature Learning on Point Sets in a Metric Space. In Guyon, I.; Luxburg, U. V.; Bengio, S.; Wallach, H.; Fergus, R.; Vishwanathan, S.; and Garnett, R., eds., *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc.
- Ravanbakhsh, S.; Schneider, J.; and Poczos, B. 2017. Deep Learning with Sets and Point Clouds. arXiv:1611.04500.
- Sabour, S.; Frosst, N.; and Hinton, G. E. 2017. Dynamic Routing Between Capsules. In Guyon, I.; Luxburg, U. V.; Bengio, S.; Wallach, H.; Fergus, R.; Vishwanathan, S.; and Garnett, R., eds., *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc.
- Skianis, K.; Nikolentzos, G.; Limnios, S.; and Vazirgiannis, M. 2020. Rep the Set: Neural Networks for Learning Set Representations. In *The 23rd International Conference on Artificial Intelligence and Statistics*, 1410–1420. PMLR.
- Sukhbaatar, S.; Grave, É.; Bojanowski, P.; and Joulin, A. 2019. Adaptive Attention Span in Transformers. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, 331–335.
- van der Maaten, L.; and Hinton, G. 2008. Visualizing Data using t-SNE. *Journal of Machine Learning Research*, 9(11).
- Wagstaff, E.; Fuchs, F.; Engelcke, M.; Posner, I.; and Osborne, M. A. 2019. On the Limitations of Representing Functions on Sets. In *Proceedings of the 36th International Conference on Machine Learning*, 6487–6494. PMLR.
- Wu, Z.; Song, S.; Khosla, A.; Yu, F.; Zhang, L.; Tang, X.; and Xiao, J. 2015. 3D ShapeNets: A Deep Representation for Volumetric Shapes. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 1912–1920.
- Yan, X.; Zheng, C.; Li, Z.; Wang, S.; and Cui, S. 2020. PointASNL: Robust Point Clouds Processing Using Nonlocal Neural Networks With Adaptive Sampling. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 5589–5598.
- Yi, L.; Kim, V. G.; Ceylan, D.; Shen, I.-C.; Yan, M.; Su, H.; Lu, C.; Huang, Q.; Sheffer, A.; and Guibas, L. 2016. A Scalable Active Framework for Region Annotation in 3D Shape Collections. *ACM Transactions on Graphics*, 35(6): 1–12.
- Zaheer, M.; Kottur, S.; Ravanbakhsh, S.; Poczos, B.; Salakhutdinov, R. R.; and Smola, A. J. 2017. Deep Sets.



In Guyon, I.; Luxburg, U. V.; Bengio, S.; Wallach, H.; Fergus, R.; Vishwanathan, S.; and Garnett, R., eds., *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc.

Zhao, H.; Jia, J.; and Koltun, V. 2020. Exploring Self-Attention for Image Recognition. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 10076–10085.

Zhao, H.; Jiang, L.; Jia, J.; Torr, P. H.; and Koltun, V. 2021. Point Transformer. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 16259–16268.