

---

# PointMask: Towards Interpretable and Bias-Resilient Point Cloud Processing

---

Saeid Asgari Taghanaki<sup>1</sup> Kaveh Hassani<sup>\*1</sup> Pradeep Kumar Jayaraman<sup>\*2</sup> Amir Hosein Khasahmadi<sup>1</sup>  
Tonya Custis<sup>3</sup>

## Abstract

Deep classifiers tend to associate a few discriminative input variables with their objective function, which in turn, may hurt their generalization capabilities. To address this, one can design systematic experiments and/or inspect the models via interpretability methods. In this paper, we investigate both of these strategies on deep models operating on point clouds. We propose *PointMask*, a model-agnostic interpretable information-bottleneck approach for attribution in point cloud models. PointMask encourages exploring the majority of variation factors in the input space while gradually converging to a general solution. More specifically, PointMask introduces a regularization term that minimizes the mutual information between the input and the latent features used to mask out irrelevant variables. We show that coupling a PointMask layer with an arbitrary model can discern the points in the input space which contribute the most to the prediction score, thereby leading to interpretability. Through designed bias experiments, we also show that thanks to its gradual masking feature, our proposed method is effective in handling data bias.

## 1. Introduction

The performance of deep neural networks is usually measured based on their predictive behavior on a validation/test set. However, evaluating a model’s performance on a single dataset can hardly capture its underlying behavior—even if the dataset is large enough (Geirhos et al., 2020). Therefore, having lucid explanations about predictions made by neural networks is critical for many applications. For example,

in the image classification task, one cannot draw a reasonable conclusion solely based on predicted class probabilities (Geirhos et al., 2020) where variations in background or textures can completely change the predictions (Beery et al., 2018; Rosenfeld et al., 2018). This critical flaw, which is mostly hidden, arises due to bias in training data (Torralba & Efros, 2011). Tasks such as image segmentation or object detection, on the other hand, where the output has a perceptible relation with the input, are inherently more explainable since one can infer if a model behaves anomalously by looking at a segmentation mask or a detected bounding box.

Deep classifiers are effective in finding a few—but not most/all—discriminative variables in the input (Geirhos et al., 2018). Relying on a few and often correlated variables can lead to poor generalization when the variables are absent at test time due to a shift in data distribution (Jo & Bengio, 2017; Geirhos et al., 2020). Uncovering biases in training data and inspecting whether a deep model converges to only a few input variables are two equally important advantages/applications of interpretability methods.

Despite the growing interest in interpreting 2D deep vision models, a scant effort has been made in interpreting deep networks processing point clouds, i.e., sparse order-invariant sets of interacting points representing 3D geometric data. Similar to images, point cloud datasets can be biased towards a specific pattern/feature. For example, a 3D sensor such as LiDAR may add a deliberate noise pattern to samples. In a multi-sensory setting, this can contribute to a *fake* improvement of classification performance when all/most samples of a specific class are collected from a sensor with a particular and consistent noise pattern which is equivalent to “context” (Beery et al., 2018) or “texture” (Geirhos et al., 2018; Baker et al., 2018) biases in images. An ideal model would learn multiple features/variables in a balanced way such that the absence of one feature at test time would not cause a drastic failure.

Interpretability methods can be categorized into three categories based on the training phase that they are applied to. Some methods are applied prior to training. Analyzing datasets for possible biases using clustering techniques are examples of such methods. Some other methods are applied during the training. Self-explanatory models which have

<sup>\*</sup>Equal contribution <sup>1</sup>Autodesk AI Lab, Toronto, Canada <sup>2</sup>Autodesk Research, Toronto, Canada <sup>3</sup>Autodesk AI Lab, San Francisco, USA. Correspondence to: Saeid Asgari <saeid.asgari.taghanaki@autodesk.com>.

interpretability modules are examples of these methods (Zhmoginov et al., 2019; Zhang et al., 2018; Fong & Vedaldi, 2017; Dabkowski & Gal, 2017). Finally, post-hoc techniques are applied to trained models (Simonyan et al., 2013; Selvaraju et al., 2017; Smilkov et al., 2017; Sundararajan et al., 2017; Springenberg et al., 2014; Schulz et al., 2020). Our work falls into the second group.

In this paper, we focus on interpreting deep models that process point clouds with the ultimate goal of adding robustness against potential dataset biases. Our approach can be coupled with various network architectures without any constraints. We opt for simplicity and adopt the commonly used PointNet (Qi et al., 2017a) architecture for our study. Inspired by InfoMask (Taghanaki et al., 2019), an information-bottleneck approach (Alemi et al., 2016) for semi-supervised object localization, we design a model which detects/visualizes input variables that PointNet relies on the most to make predictions. We show that adding a masking layer to PointNet not only provides interpretability but also increases the model’s robustness against *bias* in training data and improves the model’s performance on predicting the classes of randomly rotated objects. In summary, we make the following contributions:

- We extend InfoMask (Taghanaki et al., 2019) to interpret deep models operating on point clouds. We call the introduced model *PointMask* which learns to mask out input points with a negligible contribution to the model’s predictions while encourages for exploring the majority of the input variables.
- We also introduce *PointMap*, a variant of PointMask which instead of masking, learns to map the points into a new space, and study both models’ effectiveness for unbiased point cloud processing.
- Finally, we show that removing extra input variables introduces a regularization effect which increases the model’s robustness to random rotations.

## 2. Related Work

**Interpretability on point clouds.** Deep models introduced in the literature for processing point clouds primarily use visualizations to gain insights about the learned representations. In pioneer works such as PointNet (Qi et al., 2017a), PointNet++ (Qi et al., 2017b), and dynamic graph CNN (Wang et al., 2019), the intermediate representations are visualized by projecting them back into the point set space, whereas in (Zhao et al., 2019; Hassani & Haley, 2019), the evolution of learned representations is visualized through the training iterations. FoldingNet (Yang et al., 2018) addresses interpretability from two aspects: (1) clear geometric interpretation by imposing a *virtual force* to de-

form a 2D grid lattice onto a 3D object surface, and (2) visualizing gradual change of the folding forces.

To the best of our knowledge, C-PointNet (Zhang et al., 2019; Huang et al., 2019) is the only relevant work focusing on the interpretability of point cloud models. It aggregates the point features learned by PointNet into a class-attentive global feature and generates class-attentive response maps to explain the decision making process in PointNet. Our approach, on the other hand, introduces a differentiable layer before the encoder that learns to mask out the points with negligible contributions by maximizing mutual information between the masked points and the class labels. It is noteworthy that our proposed module can be integrated into any other point set encoder without any constraints. We opt for simplicity and adopt the commonly used PointNet architecture in this work.

**Designing interpretable models.** Unlike the high interest in post-hoc interpretability methods (Simonyan et al., 2013; Selvaraju et al., 2017; Smilkov et al., 2017; Sundararajan et al., 2017; Springenberg et al., 2014; Schulz et al., 2020), only a few works focus on adding interpretable components to deep models since it results in accuracy degradation on benchmark test sets. (Zhang et al., 2018) modified convolution layers using masks to obtain sharper feature maps. In (Zhmoginov et al., 2019), a hard attention mechanism based on mutual information is introduced that detects the most discriminative areas of the input. In the work done by (Fong & Vedaldi, 2017), a framework is proposed that learns masks to find parts of an image that influence the classifier’s decision the most. Similarly, (Dabkowski & Gal, 2017) apply a masking model to manipulate the scores of a pre-trained classifier by masking salient parts of the input. Nevertheless, these methods are all designed for RGB images whereas our method is specifically designed for point clouds. Moreover, the behavior of these models in the presence of data bias is unclear. We, on the other hand, analyze our method with different levels of input bias. Further, it is worth noting that our method does not lead to a decrease in overall classification accuracy.

## 3. Method

Given a point set  $\mathcal{P} = \{\mathcal{P}_i | i = 1, \dots, n\} \in \mathbb{R}^{n \times 3}$ , where  $\mathcal{P}_i \in \mathbb{R}^3$  denotes  $(x, y, z)$  coordinates of point  $i$ , our goal is to calculate a set of variational importance probabilities  $\mathcal{J} = \{\mathcal{J}_i | i = 1, \dots, n\} \in \mathbb{R}^n$  corresponding to  $\mathcal{P}$  by optimizing:

$$\mathcal{L}(\omega) = I(\mathcal{J}, \mathcal{Y}; \omega) - \alpha I(\mathcal{J}, \mathcal{P}; \omega) \quad (1)$$

where  $I$  is mutual information estimator,  $\omega$  is network parameters,  $\alpha$  is a scalar weight, and  $\mathcal{Y}$  is class label of  $\mathcal{P}$ . We

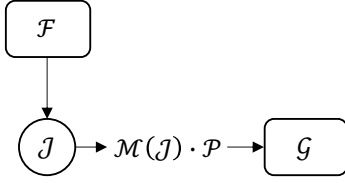


Figure 1. PointMask: the entire model is learned end-to-end.  $\mathcal{F} = p(\mathcal{J}|\mathcal{P}; \omega)$  and  $\mathcal{G} = p(\mathcal{Y}|\mathcal{M}(\mathcal{J})\cdot\mathcal{P}; \omega)$ .

rewrite  $I(\mathcal{J}, \mathcal{Y}; \omega)$  (and similarly  $I(\mathcal{J}, \mathcal{P}; \omega)$ ) as:

$$I(\mathcal{J}, \mathcal{Y}; \omega) = \int p(\mathcal{J}, \mathcal{Y}; \omega) \log \frac{p(\mathcal{J}, \mathcal{Y}; \omega)}{p(\mathcal{J}; \omega)p(\mathcal{Y}; \omega)} d\mathcal{P}d\mathcal{Y} \quad (2)$$

We assume that the underlying distribution of  $\mathcal{J}$  is a normal distribution  $\mathcal{N}(\mu_{\mathcal{J}}, \sigma_{\mathcal{J}})$  and learn the distribution parameters  $\mu \in \mathbb{R}^n$  and  $\sigma \in \mathbb{R}^n$  using function  $\mathcal{F}$  implemented as a feed-forward neural network. To sample  $\mathcal{J}$ , we use reparameterization trick, i.e.,  $\mathcal{J} = \mathcal{F}(\mathcal{P}, \epsilon) = \mu_{\mathcal{J}} + \sigma_{\mathcal{J}}\epsilon$ . We define a masking function  $\mathcal{M} = \text{ReLU}(\sigma(\mathcal{J}) - t)$  to discard probabilities less than a threshold  $t$  where  $\sigma$  is sigmoid activation and ReLU is ReLU activation with upper-bound of 1. We multiply  $\mathcal{M}$  by  $\mathcal{P}$  to remove less critical points in the input and forward the remaining points to a point set classifier  $\mathcal{G}$  which is PointNet in our case. The proposed model is shown in Figure 1.

The loss function consists of a classification loss and a regularization as follows:

$$\mathcal{L} = \frac{1}{N} \sum_{n=1}^N \mathbb{E}_{\epsilon \sim p(\epsilon)} [-\log q(\mathcal{Y}_n | \mathcal{M}(\mathcal{P}_n, \epsilon))] + \alpha \text{KL}(p(\mathcal{J}|\mathcal{P}_n) \parallel r(\mathcal{J})) \quad (3)$$

where  $N$  is the number of training samples,  $q(\cdot)$  is the variational approximation function,  $r(\mathcal{J})$  is the variational approximation to the marginal  $p(\mathcal{J}) = \int p(\mathcal{P})p(\mathcal{J}|\mathcal{P})d\mathcal{P}$ , and  $\text{KL}(\cdot, \cdot)$  is the KullbackLeibler divergence. The scalar weight  $\alpha$  controls the level of information that is passed through from the input.  $\alpha$  is critical for encouraging the model to *not* over-fit to a few input variables but to explore most/all of them leading to better generalization in a new environment where a few of the variables are absent.

Alternatively, we can change  $\mathcal{M}$  from a masking function into a mapping function by modifying the dimension of  $\mathcal{J}$  from  $\mathbb{R}^n$  to  $\mathbb{R}^{n \times 3}$ , and treating it as a per-point translation function. Therefore,  $\mathcal{M} = \text{ReLU}(\sigma(\mathcal{J}) - t)$  becomes  $\mathcal{M} = \mathcal{J}$ . This can be summed with  $\mathcal{P}$  to map it into a new space as controlled by Equation 3. We call this variant *PointMap*.

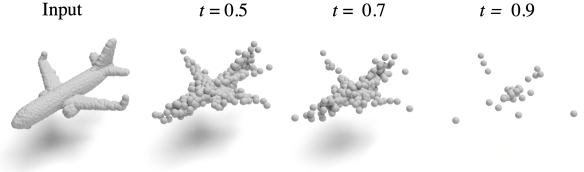


Figure 2. PointMask with different threshold values; increasing the threshold leads to detecting key elements without using any prior knowledge.

## 4. Experiments

In the following subsections, we show PointMask’s capability to interpret the classifier’s decision by detecting the important points. We then examine how it behaves given a biased training set, and finally, we test whether the masking step of PointMask helps to differentiate similar objects which are often misclassified.

We evaluate our model on ModelNet40 and ModelNet10 shape classification benchmarks (Wu et al., 2015). ModelNet40 consists of 12,311 CAD models from 40 man-made object categories, split into train and test sets of sizes 9,843 and 2,468, respectively. ModelNet10 is a subset of ModelNet40 and contains 4899 samples split into 3991 training and 908 test samples.

### 4.1. Interpreting by Detecting Important Input Variables

The first goal of this work is to visualize/detect important input variables to answer the question: “Which subset(s) of a given input point cloud contribute(s) most to the final prediction?” To this end, we propose to place a trainable variational masking layer before PointNet which learns to remove less critical points (see section 3). As shown in Figure 2, increasing the masking threshold improves the detection of the points that correlate the most with the object class.

We also compare PointMask with PointNet trained with different types of data augmentation and tested on aligned samples. Results reported in Table 4.1 suggest that PointMask in addition to providing interpretability, outperforms PointNet on the ModelNet40 dataset. We observed that augmenting the training set with random rotations results in an accuracy drop of  $\sim 1.0\%$  on both models lose (see section 4.3).

### 4.2. Unbiased Feature Learning

A deep classifier might converge to one or a few input variables/patterns such as texture, shape, etc. To study this on point clouds, we design a controlled experiment and systematically inject unique and constant patterns to each object

Table 1. ModelNet40 classification results for different models trained with data augmentation and tested on aligned samples. J, R1, R3, and A refer to jitter, rotation along a single axis, rotation along all 3 axes, and aligned samples respectively.

	Trained with - Tested on		
	J - A	JR1 - A	JR3 - A
PointNet	89.61	89.45	88.39
PointMask	89.73	89.37	88.88



Figure 3. A toy example depicting single point bias. The blue, black, and red points represent input, masked, and bias points, respectively. The first row shows samples with bias from different categories while the second and third row visualize the same samples which are randomly masked. In some samples bias and some parts of the objects have been masked together.

category. The red points placed in different coordinates in the first row of Figure 3 show the systematic patterns that are added to each class during the training. We perform a similar experiment using ModelNet10 dataset (Wu et al., 2015). We add a point cloud alphabet (Figure 5, (a)) to each category of shapes, e.g., letter “I” is concatenated with all desks where “I” is sampled with different number of points  $\{1, 50, 100, 256\}$  to increase the strength of the pattern (Figure 5 (b-e)). The goal here is to test whether PointNet attends to any part of the object or just overfits to the unique pattern. With the presence of patterns at test time, PointNet obtains  $\sim 100\%$  classification accuracy. However, as shown in Figure 4, when the patterns are removed at test time, its performance drops to random accuracy, even with a single point pattern. This indicates that regardless of the other 2048 points, PointNet only relies on a single point bias pattern.

As a simple remedy for constant bias patterns, we randomly mask  $n\%$  of the points of each object where  $n \in [10, 70]$ . The hope here is that some/all bias patterns might be removed in some training samples (see the second and third rows of Figure 3), thereby forcing a model to look at other variables, i.e., object points. As shown in Figure 4, this approach turns out to be effective for weak bias patterns,

however, when the bias level increases to 256 points, its classification accuracy on ModelNet10 decreases to 42%.

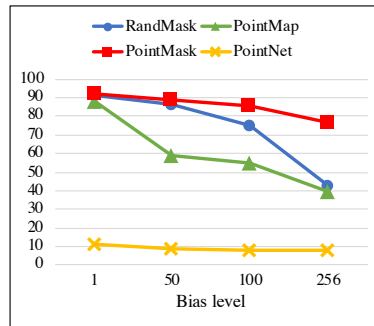


Figure 4. ModelNet10 classification results with different levels of bias added during training and removed at test time.

Next, we train PointMask and PointMap using the same biased train set and evaluate them with a bias-free test set. As shown in Figure 4, PointMask performs reasonably well despite the increase in the bias level, e.g., for the strongest bias level, it achieves an accuracy of 76% which is a 34% absolute improvement over RandMask. The intuition for this improvement is that PointMask samples masks from a normal distribution and the KL divergence term in Equation 3 enforces the masks to have a minimum amount of information about the input points; on the other hand, the classification loss enforces the sampled masks to have maximum information about class labels. This min-max game prevents the model from over-fitting to a premature solution. Therefore the model has the chance to explore a diverse set of masks per object and factors of variation during its convergence process.

We observed that for strong bias patterns, PointMap is not as effective as PointMask. This is because, in order to obliterate the bias, PointMap has to apply a large translation per point and since there is no prior knowledge on the “objectness” of the samples, this translations will morph the object to something else. Therefore, there will be no unique features across the classes left for the model to use to improve classification accuracy. This phenomenon is visualized in Figure 6 where PointMask is able to eliminate the bias by removing points, whereas for the object processed by PointMap, the bias is still present/detectable.

### 4.3. Regularization Effect on Rotated Objects

We also investigated whether translating/masking irrelevant points can lead to a better classification accuracy once models are tested on rotated objects. We noticed when PointNet is trained and tested with randomly rotated samples along all three  $\{x, y, z\}$  axes, its accuracy drops from 89.6% to 76.9% even though it is augmented with random rotations during training. On the other hand, PointMask improves performance from 76.9% to 82.2%. We speculate this is partly



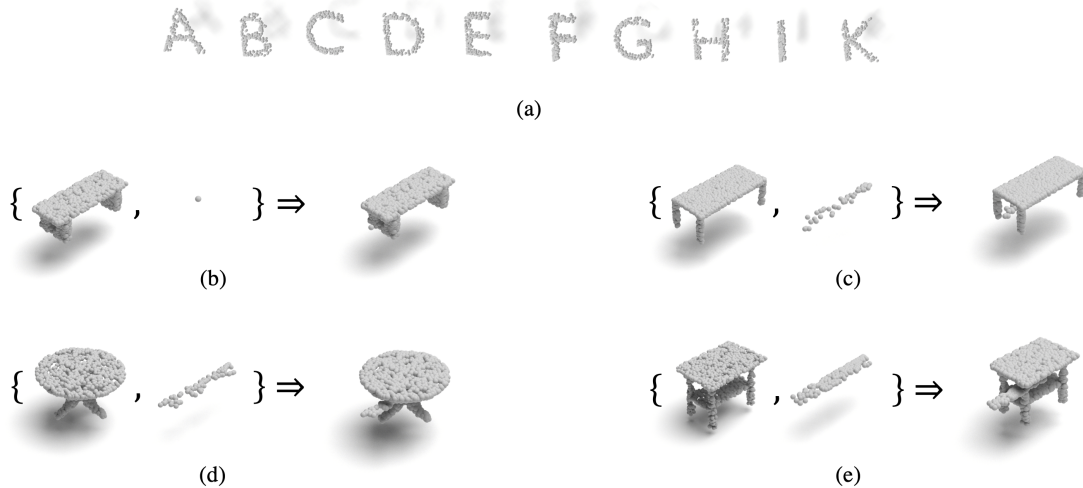


Figure 5. Generated biased samples from ModelNet10 dataset. (a) point cloud alphabet set which we add to objects as bias, (b)–(e) biased samples with 1, 50, 100, and 256 bias points, respectively.

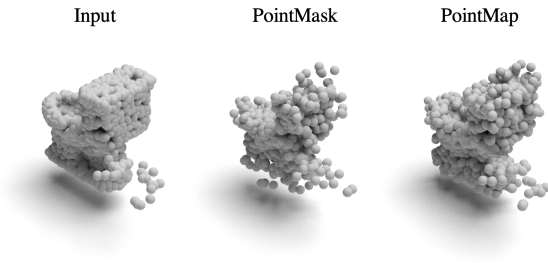


Figure 6. Visualization of processed biased samples from ModelNet10 with different methods.

because of improper regularization in PointNet, i.e., when we removed dropout layers from PointNet, the accuracy increased from 76.9% to 81.1% (see Table 4.3). Nevertheless, even with this modification, it is still outperformed by PointMask.

In addition, input rotations can confuse the model. Some object categories in ModelNet40 (e.g., bookshelf and bed) look considerably similar when viewed from similar orientations. We speculate that PointMap and PointMask help to differentiate similar objects by making them look different via point translation and removal, respectively, even if they are in the same pose. In Figure 7, we show that by increasing  $n$  in top- $n$  accuracy, all models achieve almost similar performance. However, looking at top-1 and top-2 accuracy indicates that PointNet was confused among the top few classes, whereas PointMap and PointMask obtained reasonable performance.

We further examine the per-class accuracy by different models using a confusion matrix. As shown in Table 4.3, PointMap and PointMask improve accuracy on several classes. We observe that there is a big confusion

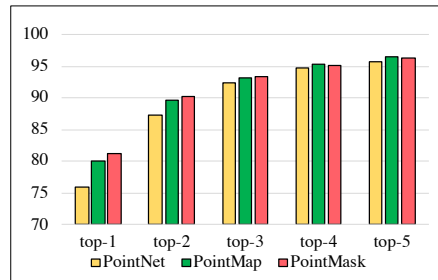


Figure 7. Top- $n$  classification accuracy on ModelNet40.

Table 2. Classification accuracy for lowest scoring ModelNet40 classes.

Class	PointNet	PointMap	PointMask
flowerpot	00.0	20.0	25.0
cup	35.0	50.0	55.0
radio	30.0	40.0	60.0
wardrobe	5.00	0.40	40.0

between flowerpot (0.0) and plant (50.0) classes in PointNet, which goes down to 40.0 in both PointMap and PointMask. Similarly, the confusion of the cup (35.0) class is highest with the vase (45.0) class, and goes down to 25.0 and 30.0 in PointMap and PointMask, respectively. We observe similar trends in other poorly performing classes as well.

In Figure 8, we visualize the output of the input transformation module in PointNet which consists of the T-Net and rotation operation. As shown, different rotations are obtained when the same input is processed by different methods.

Table 3. ModelNet40 classification results for different models trained and tested with randomly rotated objects. Rotation here can be about any random 3D axis.

	Original	Without dropout
PointNet	76.87	81.13
RandMask	75.20	80.64
PointMap	81.05	NA
PointMask	82.18	NA

#### 4.4. Implementation Details

For all models, we set the number of points to 2048 and their dimension to 3, i.e., Euclidean coordinates. We optimized the models for 500 epochs using Adam (Kingma & Ba, 2014) with learning rate set to 0.0001, and mini-batch size of 32. We provide pseudocode for PointMask in Listing 1.

```

1 def mask_relu(x, threshold = 0.5):
2     x = sigmoid(x)
3     x = relu(x - threshold)
4     x = clamp(x, min=0, max=1)
5     return x
6
7 def kl_divergence(mu, log_var):
8     tmp = 1.0 + log_var - mu * mu - exp(
9         log_var)
10    kl_batch = -0.5 * sum(tmp, axis=-1)
11    return alpha * mean(kl_batch)
12
13 # PointMask
14 num_points = 2048
15 feature_dim = 3 # (x, y, z) coordinates
16 P = ... # pointcloud of shape (num_points,
17         feature_dim)
18 l = Conv1DReLU(channels=64)(P)
19 l = Conv1DReLU(channels=128)(l)
20 l = Conv1DReLU(channels=1024)(l)
21 l = MaxPool1D(pool_size=2048)(l)
22
23 # Reparameterization trick
24 mu = Conv1D(channels=2048)(l)
25 log_var = Conv1D(channels=2048)(l)
26 loss_kl = kl_divergence(mu, log_var)
27 sigma = exp(0.5 * log_var)
28 eps = random_normal(stddev=1.0, shape=(
29     num_points, 1, 2048))
30 eps = sigma * eps
31 J = mu + eps
32
33 # Masking input points
34 M = mask_relu(J)
35 masked_input = P * M
36
37 # Forward masked points to PointNet
38 pred = G(masked_input) # G = PointNet
39 loss_ce = cross_entropy(pred, ground_truth)
40 L = loss_kl + loss_ce

```

Listing 1. Pseudocode for PointMask.

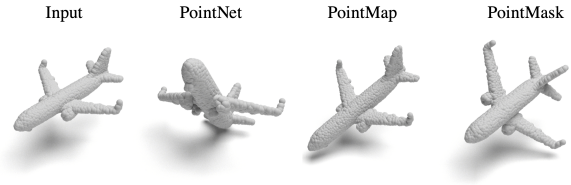


Figure 8. The output of input transformation module (T-Net and rotation) in PointNet for different methods.

## 5. Conclusion

We introduced PointMask, an information-bottleneck interpretability layer that can be integrated into any point cloud classification model. We showed that PointMask not only provides interpretability but also brings robustness against deliberate bias patterns which might not be perceptible to humans. We also demonstrated that the regularization effect of the proposed model is more effective compared to dropout. Finally, we showed that our method reduces confusion among classes that might be similar when arbitrarily rotated. As a future direction, we are planning to apply the proposed method to key point/landmark detection in point clouds.

## Acknowledgements

We would like to thank Aditya Sanghi and Ara Danielyan for their useful input through the course of this research.

## References

- Alemi, A. A., Fischer, I., Dillon, J. V., and Murphy, K. Deep variational information bottleneck. *arXiv preprint arXiv:1612.00410*, 2016.
- Baker, N., Lu, H., Erlikhman, G., and Kellman, P. J. Deep convolutional networks do not classify based on global object shape. *PLoS computational biology*, 14(12): e1006613, 2018.
- Beery, S., Van Horn, G., and Perona, P. Recognition in terra incognita. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pp. 456–473, 2018.
- Dabkowski, P. and Gal, Y. Real time image saliency for black box classifiers. In *Advances in Neural Information Processing Systems*, pp. 6967–6976, 2017.
- Fong, R. C. and Vedaldi, A. Interpretable explanations of black boxes by meaningful perturbation. In *Proceedings of the IEEE International Conference on Computer Vision*, pp. 3429–3437, 2017.
- Geirhos, R., Rubisch, P., Michaelis, C., Bethge, M., Wichmann, F. A., and Brendel, W. Imagenet-trained cnns are

- biased towards texture; increasing shape bias improves accuracy and robustness. *arXiv preprint arXiv:1811.12231*, 2018.
- Geirhos, R., Jacobsen, J.-H., Michaelis, C., Zemel, R., Brendel, W., Bethge, M., and Wichmann, F. A. Shortcut learning in deep neural networks. *arXiv preprint arXiv:2004.07780*, 2020.
- Hassani, K. and Haley, M. Unsupervised multi-task feature learning on point clouds. In *Proceedings of the IEEE International Conference on Computer Vision*, pp. 8160–8171, 2019.
- Huang, S., Zhang, B., Shen, W., and Wei, Z. A claim approach to understanding the pointnet. In *Proceedings of the 2019 2nd International Conference on Algorithms, Computing and Artificial Intelligence*, pp. 97–103, 2019.
- Jo, J. and Bengio, Y. Measuring the tendency of cnns to learn surface statistical regularities. *arXiv preprint arXiv:1711.11561*, 2017.
- Kingma, D. P. and Ba, J. Adam: A method for stochastic optimization, 2014.
- Qi, C. R., Su, H., Mo, K., and Guibas, L. J. Pointnet: Deep learning on point sets for 3d classification and segmentation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 652–660, 2017a.
- Qi, C. R., Yi, L., Su, H., and Guibas, L. J. Pointnet++: Deep hierarchical feature learning on point sets in a metric space. In *Advances in neural information processing systems*, pp. 5099–5108, 2017b.
- Rosenfeld, A., Zemel, R., and Tsotsos, J. K. The elephant in the room. *arXiv preprint arXiv:1808.03305*, 2018.
- Schulz, K., Sixt, L., Tombari, F., and Landgraf, T. Restricting the flow: Information bottlenecks for attribution. *arXiv preprint arXiv:2001.00396*, 2020.
- Selvaraju, R. R., Cogswell, M., Das, A., Vedantam, R., Parikh, D., and Batra, D. Grad-cam: Visual explanations from deep networks via gradient-based localization. In *Proceedings of the IEEE international conference on computer vision*, pp. 618–626, 2017.
- Simonyan, K., Vedaldi, A., and Zisserman, A. Deep inside convolutional networks: Visualising image classification models and saliency maps. *arXiv preprint arXiv:1312.6034*, 2013.
- Smilkov, D., Thorat, N., Kim, B., Viégas, F., and Wattenberg, M. Smoothgrad: removing noise by adding noise. *arXiv preprint arXiv:1706.03825*, 2017.
- Springenberg, J. T., Dosovitskiy, A., Brox, T., and Riedmiller, M. Striving for simplicity: The all convolutional net. *arXiv preprint arXiv:1412.6806*, 2014.
- Sundararajan, M., Taly, A., and Yan, Q. Axiomatic attribution for deep networks. In *Proceedings of the 34th International Conference on Machine Learning-Volume 70*, pp. 3319–3328. JMLR. org, 2017.
- Taghanaki, S. A., Havaei, M., Berthier, T., Dutil, F., Di Jorio, L., Hamarneh, G., and Bengio, Y. Infomask: Masked variational latent representation to localize chest disease. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pp. 739–747. Springer, 2019.
- Torralba, A. and Efros, A. A. Unbiased look at dataset bias. In *CVPR 2011*, pp. 1521–1528. IEEE, 2011.
- Wang, Y., Sun, Y., Liu, Z., Sarma, S. E., Bronstein, M. M., and Solomon, J. M. Dynamic graph cnn for learning on point clouds. *ACM Transactions on Graphics (TOG)*, 38(5):1–12, 2019.
- Wu, Z., Song, S., Khosla, A., Yu, F., Zhang, L., Tang, X., and Xiao, J. 3d shapenets: A deep representation for volumetric shapes. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 1912–1920, 2015.
- Yang, Y., Feng, C., Shen, Y., and Tian, D. Foldingnet: Point cloud auto-encoder via deep grid deformation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 206–215, 2018.
- Zhang, B., Huang, S., Shen, W., and Wei, Z. Explaining the pointnet: What has been learned inside the pointnet? In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*, pp. 71–74, 2019.
- Zhang, Q., Nian Wu, Y., and Zhu, S.-C. Interpretable convolutional neural networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 8827–8836, 2018.
- Zhao, Y., Birdal, T., Deng, H., and Tombari, F. 3d point capsule networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 1009–1018, 2019.
- Zhmoginov, A., Fischer, I., and Sandler, M. Information-bottleneck approach to salient region discovery. *arXiv preprint arXiv:1907.09578*, 2019.