

Metric Learning for 3D Point Clouds Using Optimal Transport

Siddharth Katageri

Srinjay Sarkar

Charu Sharma

Machine Learning Lab, International Institute of Information Technology, Hyderabad, India

siddharth.katageri@research.iiit.ac.in, srinjay95@gmail.com, charu.sharma@iiit.ac.in

Abstract

Learning embeddings of any data largely depends on the ability of the target space to capture semantic relations. The widely used Euclidean space, where embeddings are represented as point vectors, is known to be lacking in its potential to exploit complex structures and relations. Contrary to standard Euclidean embeddings, in this work, we embed point clouds as discrete probability distributions in Wasserstein space. We build a contrastive learning setup to learn Wasserstein embeddings that can be used as a pre-training method with or without supervision towards any downstream task. We show that the features captured by Wasserstein embeddings are better in preserving the point cloud geometry, including both global and local information, thus resulting in improved quality embeddings. We perform exhaustive experiments and demonstrate the effectiveness of our method for point cloud classification, transfer learning, segmentation, and interpolation tasks over multiple datasets including synthetic and real-world objects. We also compare against recent methods that use Wasserstein space and show that our method outperforms them in all downstream tasks. Additionally, our study reveals a promising interpretation of capturing critical points of point clouds that makes our proposed method self-explainable.

1. Introduction

Recent years have seen major advancements in 3D point cloud representation learning. It has gained prominence in a wide spectrum of areas such as robotics [9], computer vision [19], and animation [13] with a broad range of applications including shape synthesis and modeling [25], autonomous driving [8], and indoor navigation [28].

Metric learning for good quality point cloud embeddings is a crucial problem given the unique set of challenges associated with 3D data, from processing point clouds in various forms to learning in different spaces. Processing and developing learning methods for point clouds is one of the major challenges due to their irregular, unstructured and un-

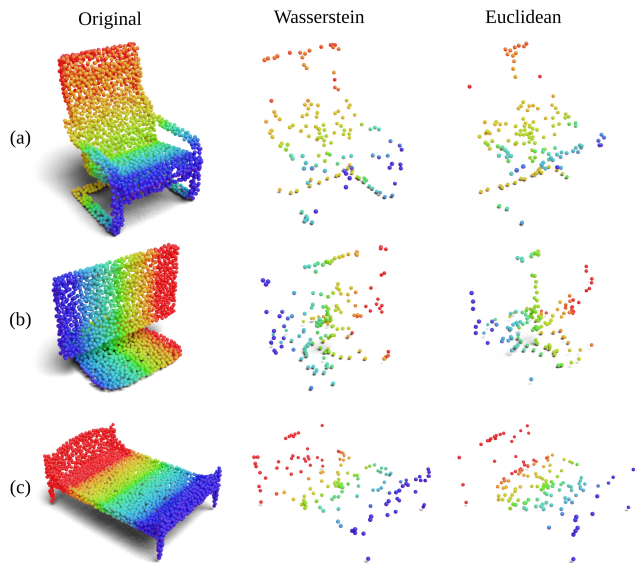


Figure 1. **Critical Points** contribute to the point cloud embedding by capturing global geometry. (a), (b) and (c) represent original point cloud (first column), critical point set for Wasserstein space (second column) and Euclidean space (third column) for three examples (Chair, Monitor & Bed).

ordered nature.

Earlier methods process point clouds by converting them into regular structures, like volumetric representations [9, 23] or 2D image projections [16, 19] to employ well explored powerful convolutional techniques. However, these transformations either incur loss of information or require high memory and computational complexity. Later, methods have been developed to learn representations by directly using raw point clouds [15, 17, 22]. These methods either process each point individually or try to infer features from local regions in a point cloud.

The common choice of recent 3D point cloud representation learning methods is to operate and represent point clouds as point vectors in *Euclidean spaces*, where relation between data points is depicted by either angle or distance. We all know that the embedding space largely determines the quality of embeddings, as it depends on how well the

target space can capture the structure of data. Euclidean space is confined in its potential to capture *complex structure* and possible semantic relations. Realizing these drawbacks, many works use hyperbolic space [11, 12] to capture this uncertainty and asymmetric relationship for word and graph embeddings.

As Euclidean space is constrained in its ability to represent data structures, we need to go *beyond Euclidean space* to get more expressive embeddings for point clouds. Recent studies show that many spaces can be embedded into *Wasserstein space* with low distortion [4], this reflects how flexible Wasserstein spaces are. Recently, [3] tries to mimic Wasserstein distance in Euclidean space for image embeddings to build efficient methods along with availing the *flexibility* of Wasserstein space. Also, there are some latest methods for getting point cloud embeddings using *Optimal Transport (OT) based distances*. [6], motivated by [3], proposes a method to approximate Wasserstein distance by Euclidean norm between two point cloud embeddings. Since Euclidean space is known for its limited ability, finding *isometric low-distortion* point cloud embeddings is tough. Another work by [10] presents how Optimal Transport based distances for point cloud reconstruction affect the quality of learnt embeddings. However, this method utilizes OT based distances only as a reconstruction loss, which is not enough to learn *complex shapes* and fails to capture *fine details* of point clouds.

Motivated by aforementioned limitations and inspired by [4], in this paper, we advocate for mapping point cloud as a *discrete distribution* in Wasserstein space. We build a contrastive learning setup to learn point cloud embeddings. Leveraging the idea of contrasting point clouds against each other, we intend to learn *common and distinctive features* between same and different distributions, respectively. It can be applied to both *supervised* and *self-supervised* settings. We use *Sliced Wasserstein (SW) distance* which is a low-cost approximation of Wasserstein distance due to its high computational complexity. Along with comparisons with commonly used distance measures such as L^2 norm and Cosine similarity, we also compare our method against recent works on point clouds that use OT. The main motive of this study is to examine the advantages of Wasserstein space over widely used Euclidean space for point cloud data. We restrict ourselves from comparing against other pre-training methods based on powerful transformer architectures [5, 18, 26, 27] and instead compare with works that explore ways of extracting point cloud embeddings using Wasserstein metric. We also show that our learnt features capture the point cloud structure better than Euclidean embeddings and consistently outperform all baselines in multiple 3D analysis and synthesis tasks. We argue that our approach of incorporating OT metric in a contrastive learning setup captures the underlying geometry and global shape

pertaining to critical points (as shown in Figure 1) and fine details of a point cloud.

Our Contributions:

1. To the best of our knowledge, we are the first to propose the use of OT metric as a distance measure in contrastive learning for point cloud data. Unlike, Euclidean embeddings, we represent a point cloud as a discrete distribution in the embedding space.
2. Using this representation, we design a framework for learning Wasserstein embeddings for 3D point clouds endowed by contrastive learning with Sliced Wasserstein distances.
3. We perform exhaustive experiments on a wide variety of downstream tasks over four popular datasets to validate the effectiveness of these embeddings over Euclidean ones and other baselines. We also illustrate the efficacy of Wasserstein embeddings by visualizing the 3D features captured by the model.

2. Preliminaries

In this section, we briefly present the optimal transport metric, variants of Wasserstein distance, and contrastive learning setup which are used in our proposed method.

2.1. Optimal Transport and Wasserstein Distance

Optimal transport aims to solve for the most efficient way to transport mass between two probability distributions. Formally, given two probability distributions μ and ν on a metric space \mathcal{X} , for $p \geq 1$, the p -Wasserstein distance is given by

$$W_p(\mu, \nu) = \left(\inf_{\pi \in \Pi(\mu, \nu)} \int_{\mathcal{X} \times \mathcal{X}} c(x, y)^p d\pi(x, y) \right)^{1/p} \quad (1)$$

where, π is a transport plan that defines a flow between mass from μ to locations in ν , $\Pi(\mu, \nu)$ is the joint probability distribution with the marginals μ and ν and $c(x, y)$ is the ground metric which assigns a cost of moving a unit of mass $x \in \mathcal{X}$ from μ to some location $y \in \mathcal{X}$ in ν . The cost of moving the mass in μ to match in ν according to the optimal transport plan π^* , is called the Wasserstein distance between the two distributions [21].

The above equation can also be written for discrete distributions, say $\hat{\mu} = \sum_{i=1}^m a_i \delta(x_i)$ and $\hat{\nu} = \sum_{j=1}^n b_j \delta(y_j)$ are two discrete distributions, where, $\{a_i\}; i = 1 \dots m$ and $\{b_j\}; j = 1 \dots n$ are the probability masses that should sum to 1, δ is the Dirac delta function and $\{x_i\}; i = 1 \dots m$ and $\{y_j\}; j = 1 \dots n$ are the support points in \mathbb{R}^d with m and n being the number of points in each measure. Then, the

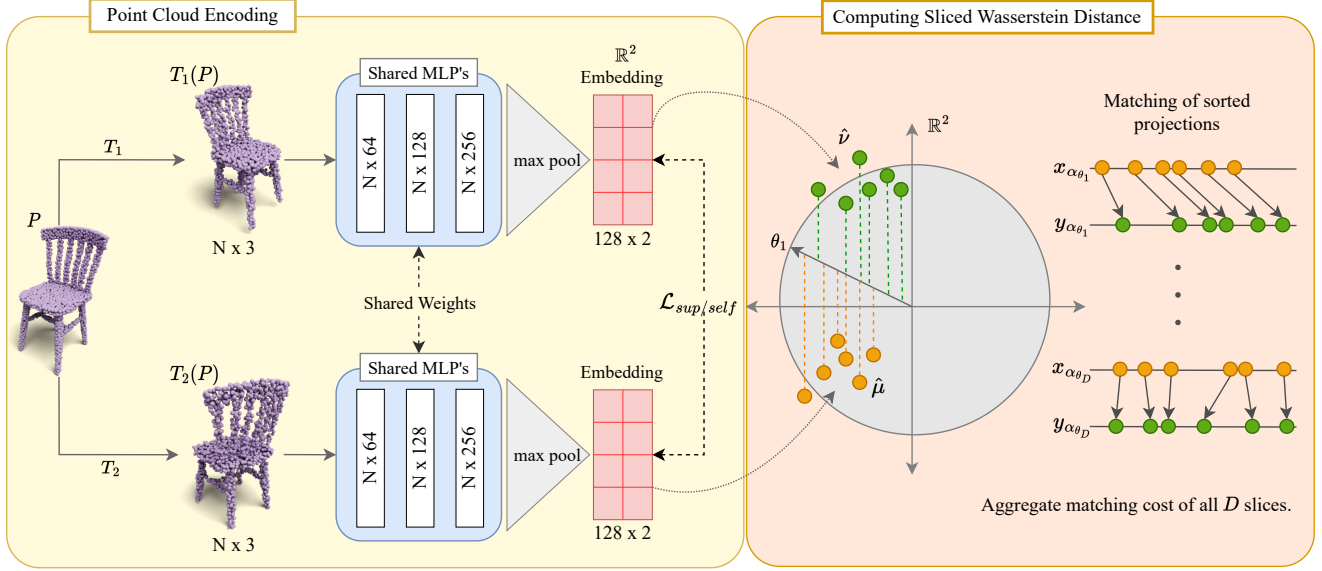


Figure 2. **Overview** of our proposed method. The two main parts are, point cloud encoding as discrete distribution (left) and computation of Sliced Wasserstein distance (right). The ground metric space is \mathbb{R}^2 . $T_1(P)$ and $T_2(P)$ are two instances of P after random transformations.

discrete version of Equation 1 is

$$W_p(\hat{\mu}, \hat{\nu}) = \left(\min_{P \in U(a,b)} \langle C^p, P \rangle \right)^{1/p} \quad (2)$$

where, $\langle \cdot, \cdot \rangle$ denotes the Frobenius dot-product, $C \in \mathbb{R}_+^{m \times n}$ is the pairwise ground metric distance, P is the coupling matrix and U is the set of all possible valid coupling matrices, i.e. $U(a, b) = \{P \in \mathbb{R}^{m \times n} : P\mathbf{1}_n = a, P^\top \mathbf{1}_m = b\}$.

Interestingly, there exists a closed-form solution for Wasserstein distance only when the distributions are one-dimensional measures with L^p norm as the cost function. The closed-form for Wasserstein distance in 1-D is [14]

$$W_p(\mu, \nu) = \left(\int_0^1 |F_\mu^{-1}(t) - F_\nu^{-1}(t)|^p dt \right)^{1/p} \quad (3)$$

where, F_μ^{-1} and F_ν^{-1} are the inverse cumulative distribution functions of μ and ν .

Generally, we are more interested in dimensions greater than one. Thus, we cannot use this closed-form solution directly to solve the OT problem efficiently. Instead, the Wasserstein distance between two measures on \mathbb{R}^d can be approximated by aggregating the 1-D Wasserstein distance between their projections over multiple directions on a unit sphere, which is called the Sliced Wasserstein distance [14]:

$$SW_p(\mu, \nu) = \left(\int_{S^{d-1}} W_p(P_{\theta, \#} \mu, P_{\theta, \#} \nu)^p d\theta \right)^{1/p} \quad (4)$$

where, $S^{d-1} = \{\theta \in \mathbb{R}^d : \|\theta\| = 1\}$ is the d -dimensional unit sphere and $P_\theta : \mathbb{R}^d \rightarrow \mathbb{R}$ is the projection. Since the

projections are now 1-D measures, we can use the closed-form solution given by Equation 3. When $m = n$, the Sliced Wasserstein distance can be easily computed by simply sorting points in 1-D measures and can be given by:

$$SW_p(\hat{\mu}, \hat{\nu}) = \left(\frac{1}{D} \sum_{k=1}^D \sum_{i=1}^m |x_{\alpha_{\theta_k}(i)} - y_{\beta_{\theta_k}(i)}|^p \right)^{1/p} \quad (5)$$

where, α_{θ_k} and β_{θ_k} are the permutation ordering in the increasing order of the support points projected to the direction θ_k with D being the total number of directions.

2.2. Contrastive Learning

Contrastive learning aims to learn an embedding space that encourages augmentations of the same input sample to have similar representations and of different samples to be dissimilar. [2] is an early example of using contrastive learning in a supervised learning setup which takes pair of samples as input to the network.

On the other hand, the contrastive loss introduced by [1] is named as SimCLR. It follows batch-wise training and is operated in self-supervised setting. For this setup, the distance is reduced between the sample and its augmentations. Later, [7] proposed the extension of SimCLR for supervised setup. It additionally aims at reducing the distance between a sample and other samples from same class in a supervised setting.

3. Our Method

In this section, we discuss our method of computing Wasserstein embeddings for point clouds in a contrastive

learning setup as shown in Figure 2. We build an in-batch contrastive learning setup which can either be fully supervised or self-supervised and can be used as a pre-training methodology for any downstream task. The goal is to represent samples from same class closer than the samples from different classes in the embedding space (larger inter-cluster and smaller intra-cluster distance). Here, the choice of embedding space plays a key role for desirable performance, as individual metric spaces can embed data differently and represent different types of semantic structure.

3.1. Contrastive Learning with Optimal Transport

Let $O = \{(P_m, l_m)\}$; $m = 1 \dots M$ be a collection of point clouds $P_m = \{p_i\}$; $i = 1 \dots N_m$, where, $p_i \in \mathbb{R}^3$ with their corresponding class labels $l_m \in L$, where $L = \{1, \dots, C\}$ is a set of class labels. Each point cloud P_m contains N_m number of points defined by 3D space points in x , y and z direction. For defining the batch-wise contrastive loss, we first randomly draw K samples from the collection O , that form a batch $B = \{(P_m, l_m)_k\}$; $k = 1 \dots K$. For every point cloud $P_m \in B$, we apply fixed set of random transformations T_1 and T_2 to get two instances of P_m (as shown in Figure 2), giving an augmented batch $B' = \{(P'_m, l_m)_{k'}\}$; $k' = 1 \dots 2K$. The augmented batch is twice the size of the original batch. The point clouds P'_m indexed at k' and $k' + 1$ are augmented version of the point cloud P_m indexed at k . As these are augmented versions of $P_{m[k]}$, their class labels are $l_{m[k']} = l_{m[k'+1]} = l_{m[k]}$.

The input to the encoder is an augmented batch B' , from which all P'_m needs to be mapped to the embedding space depending on its geometric features and appearance, with samples having same class label being closer. The encoder represents function $f : \mathbb{R}^{N_m \times 3} \rightarrow \mathcal{W}(\mathcal{X})$, that maps a point cloud P'_m to the Wasserstein space $\mathcal{W}(\mathcal{X})$, with W_p being the distance metric on $\mathcal{W}(\mathcal{X})$ and \mathcal{X} being the ground metric space. We choose \mathbb{R}^2 , \mathbb{R}^4 and \mathbb{R}^8 to be our ground metric spaces, in which the corresponding embedding z'_m of P'_m is represented as discrete distribution $\{\frac{1}{S} \cdot x_i\}$; $i = 1 \dots S$ supported by $x_i \in \mathcal{X}$ with a total of S support points, all with uniform probability mass $\frac{1}{S}$ for simplicity. In our implementation, we obtain z'_m by reshaping the encoder's output, to obtain the discrete distribution for different ground metric spaces.

Generally, the computation for exact solution of W_p is costly. To make the computation of optimal transport more tractable, we replace the distance metric W_p on Wasserstein space $\mathcal{W}(\mathcal{X})$ by the Sliced Wasserstein distance metric SW_p . SW_p is a low-cost approximation of Wasserstein distance with computational complexity being $\mathcal{O}(S \log S)$. For all our experiments, we set the value of $p = 2$ and number of slices $D = 300$.

Supervised Contrastive Loss. In the supervised setting, for any $P'_m \in B'$ indexed at k' with corresponding label

$l_{m[k']}$, the positive set is defined as $A = \{P'_m \in B' : P'_m = l_{m[k']}\}$. We define our supervised contrastive loss for learning point cloud Wasserstein embeddings as:

$$\mathcal{L}_{sup} = - \sum_{i=1}^{2K} \log \left(\frac{\sum_{\substack{j \in A \\ j \neq i}} \exp(-SW_2^2(z_i, z_j))}{\sum_{t \neq i} \exp(-SW_2^2(z_i, z_t))} \right) \quad (6)$$

The loss tries to minimize the Sliced Wasserstein distance between the embeddings represented as discrete distribution of an anchor and all the samples having the same class in the augmented batch. This can also be easily converted to a self-supervised version by making necessary modifications.

Self-Supervised Contrastive Loss. Contrary to the supervised setting, in self-supervised setting, the class label of point clouds cannot be used in any way to train the encoder. Here, the positive set of any $P'_m \in B'$ contains only the other augmentation of P'_m . If $i \in \{1 \dots 2K\}$ be the index of any $P'_m \in B'$, then, let $j(i)$ be the index of its other augmented sample. We define our self-supervised loss for learning point cloud Wasserstein embeddings as:

$$\mathcal{L}_{self} = - \sum_{i=1}^{2K} \log \left(\frac{\exp(-SW_2^2(z_i, z_{j(i)}))}{\sum_{t \neq i} \exp(-SW_2^2(z_i, z_t))} \right) \quad (7)$$

Here, only the Sliced Wasserstein distance between embeddings of an anchor and its augmented sample is minimized. Other than the augmented sample, the samples having the same class in the augmented batch are treated as negatives, which might hinder the overall optimization process depending on the batchsize.

4. Experiments

Representation that is able to capture good geometric information in a smooth latent space is generally better in various shape understanding and synthesis tasks. To demonstrate the representation power of the learned Wasserstein embeddings compared to Euclidean embeddings and other baselines, in this section, we present qualitative and quantitative evaluations on multiple tasks: point cloud classification, transfer learning, point cloud segmentation, and point cloud interpolation with both supervised and self-supervised pre-training settings.

Datasets. We use ModelNet10 (MN10) and ModelNet40 (MN40) [23] to perform experiments on classification. MN40 consists of 12311 CAD models with a total of 40 categories, where 9843 objects are used for training and 2468 for testing. We use the data provided by [17], from which we randomly sample 2048 points for each point cloud. MN10 is a subset of MN40 dataset having 10 categories. To evaluate how the learned embeddings perform on real-world data, we also conduct experiments on ScanOb-

Table 1. Results of 3D object classification with supervised and self-supervised pre-training on ModelNet10, ModelNet40 and ScanObjectNN (referred as ScanObject) datasets. WPCE and SSW-AE are unsupervised methods and cannot be evaluated for supervised pre-training (represented by “-”). Bold represents the best result.

Method	Supervised Pre-training			Self-Supervised Pre-training		
	ModelNet10	ModelNet40	ScanObject	ModelNet10	ModelNet40	ScanObject
CL+ L^2	90.85	84.64	62.82	90.63	84.72	63.51
CL+Cosine	85.90	70.42	56.11	85.90	72.64	56.45
WPCE	-	-	-	89.97	78.84	52.83
SSW-AE	-	-	-	88.88	76.86	51.29
CL+ SW_2 (Ours)	91.85	85.90	63.16	91.96	85.73	63.85

jectNN (SO) [20]. It contains object scans with partial occlusions and background, making it a challenging dataset. It has 2304 objects for training and 567 for testing from 15 categories. For part segmentation, we use ShapeNetPart (SN) [25] that consists of 16681 point clouds from 16 categories and 50 part categories in total.

Pre-training. We use a 3-layer MLP followed by a max-pooling layer as our encoder for classification and segmentation tasks. For interpolation, we consider the encoder and decoder proposed by FoldingNet [24]. In order to perform any downstream task on a particular dataset, the encoder is first pre-trained on the dataset using the contrastive loss explained in Section 3.1 with different distance metrics, followed by testing and evaluation of the desired task. Throughout the experiments, we refer the encoder trained using our method as CL+ SW_2 . We train CL+ SW_2 with ground metric dimensions 2, 4, and 8 and report the results of the best-performing network. For the transformations required in contrastive loss, intended towards forming augmented instances, we sequentially compose random scaling, rotation and point jittering. In the case of Euclidean distance metrics, the encoder function $f : \mathbb{R}^{N_m \times 3} \rightarrow \mathbb{R}^d$ maps a point cloud to d -dimensional space. We then apply l^2 -distance or cosine similarity as distance measures on these d -dimensional embeddings. For the cosine similarity metric, we remove the negative sign in the numerator for Eqs. 6 and 7. Also, note that when training the encoder with cosine similarity, the embeddings are normalized. We provided full details in the supplementary material.

Baselines. For computing Euclidean embeddings, we consider L^2 -distance and Cosine similarity as distance measures. For these Euclidean metric baselines, we train the encoder using our loss (Eqs. 6, 7) by replacing $SW_2^2(\cdot, \cdot)$ with the respective metrics. We also consider WPCE [6] and SSW-AE [10] as our baselines as they are recent techniques that use Wasserstein metric for learning point cloud embeddings and are the closest comparable works to our approach. WPCE embeds Wasserstein space into Euclidean space using a Siamese network. The network is trained in

such a way that the Euclidean distance mimics the Wasserstein distance between two point clouds. SSW-AE proposed to use SW distance and its variants (max SW and adaptive SW) for reconstruction to learn point cloud embeddings. It examines the effect of using different reconstruction metrics and losses for training an auto-encoder architecture on the learnt embeddings. For a fair comparison, all reported results (apart from interpolation) of our method and baselines are using the same encoder architecture.

Hyperparameters. We perform all our experiments on NVIDIA RTX-2080Ti GPUs using the PyTorch framework for implementing our models. We set the batch size to 16, and the learning rate as 0.001 with a step learning rate scheduler, where the learning rate is scaled by 0.7 after every 20 epoch. We use the Adam optimizer and set weight decay to 0.0001 and momentum to 0.9.

Table 2. Results of 3D object classification with supervised and self-supervised pre-training for *transfer learning* setup. WPCE and SSW-AE are unsupervised methods and cannot be evaluated for supervised pre-training (represented by “-”). Bold represents the best result.

Method	Supervised Pre-training		Self-Supervised Pre-training	
	MN10 to MN40	SN to MN40	MN10 to MN40	SN to MN40
CL+ L^2	85.37	85.81	84.27	83.83
CL+Cosine	74.51	69.12	75.32	71.47
WPCE	-	-	77.51	78.03
SSW-AE	-	-	76.05	76.66
CL+ SW_2 (Ours)	86.18	86.18	85.61	85.77

4.1. 3D Object Classification

We extract point cloud embeddings from a pre-trained encoder and use a linear SVM as our classifier for simplicity. Particularly, we fit a linear SVM classifier on the embeddings acquired by an encoder on the train split and report the overall classification accuracy on the test split. In Figure 1, we can see that features captured by Wasserstein embeddings effectively summarize the overall object geom-

etry compared to embeddings learned in Euclidean space. This property also reflects in the classification performance shown in Table 1. We can observe that for both supervised and self-supervised settings, the classification accuracy with embeddings extracted by the encoder trained with $CL+SW_2$ is higher than that of $CL+L^2$ and $CL+Cosine$. Thus, compared to Euclidean space, the performance of SW_2 is consistently better on all the datasets, which implies that embeddings learnt in Wasserstein space can increase classification accuracy.

We also show that our method is more effective compared to WPCE and SSW-AE. This improvement can be explained by the difference in the approach of extracting Wasserstein embeddings, where in, our methodology introduces usage of OT metric to directly operate in embedding space endowed by contrastive learning. It helps in learning better representations by exploiting the similarities between distributions along with utilizing the flexibility of the target Wasserstein space.

4.2. Transfer Learning

We examine the generalizing ability of the embeddings acquired by encoders trained with different distance metrics to unseen classes by performing transfer learning for point cloud classification. We follow the same process as explained in Section 4.1 for reporting the overall classification accuracy. The quantitative comparisons of transfer learning are shown in Table 2. We perform evaluation in two transfer learning settings, MN10 to MN40 and SN to MN40. Here, the encoder is pre-trained on MN10 and SN, followed by evaluation on MN40. In both settings, the model generalizes to new unseen classes by wielding the knowledge of geometry learned during training. We can see that $CL+SW_2$ consistently performs better than other distance measures and methods in both the transfer learning settings with and without supervision. Results imply that Wasserstein embeddings are better in transferring the knowledge of capturing geometry for yielding good classification performance.

4.3. 3D Object Part Segmentation

We train a 3-layer MLP network to predict a class label for all points in a point cloud, where the input to this network is the embedding provided by a pre-trained encoder. In particular, part segmentation requires a fine-grain understanding of the local geometry of the objects. Along with the global embedding of the point cloud, per-point embeddings acquired before max-pooling are stacked together and passed to the segmentation network. Note that only the segmentation network weights are optimized using the standard cross-entropy loss, and the encoder’s weights are frozen. We evaluate the performance using mIoU metric. For mIoU of each class, the IoUs of all parts from that class are averaged. Instance average mIoU is calculated by taking the

mean of IoUs for all the instances. The comparison of average instance mIoU and per class average mIoU for both supervised and self-supervised learning settings are shown in Table 3. We can see that our results outperform other distance measures and methods, implying that Wasserstein embeddings are able to capture better fine-grain local information required for the task.

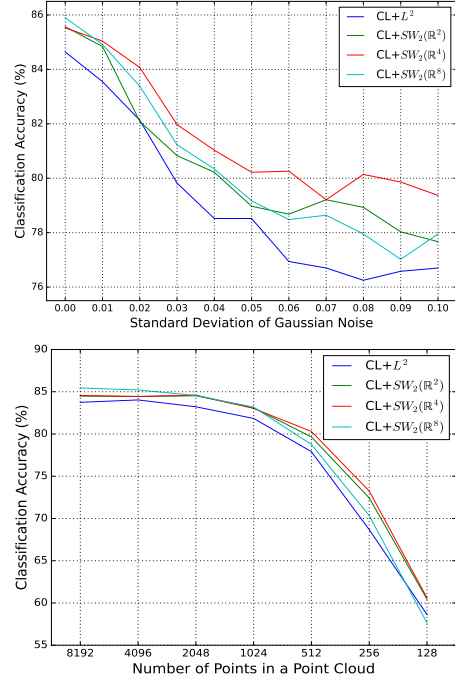


Figure 3. **Ablation Study:** Classification for noise model with Gaussian Noise (top) and points removal using random sampling (bottom) for our method $CL+SW_2$ and baseline $CL+L^2$.

4.4. Ablation Study: Classification Task

We perform point perturbation and point density variation to test their effects on the encoders pre-trained with different distance metrics and report the classification accuracy on Modelnet40 as shown in Figure 3. For the point perturbation test, we add Gaussian noise to input point clouds, with standard deviation of noise varying from 0.01 to 0.1. We can observe that for all noise levels, even with severe distortion, $CL+SW_2$ performs well than that of $CL+L^2$. This implies that discrete representation learnt in Wasserstein space is less prone to performance degradation due to noise in inputs. Further, for varying density test, we randomly sample 8192, 4096, 2048, 1024, 512, 256, and 128 points from input point clouds and perform evaluation on them. We can observe that $CL+SW_2$ consistently does better than $CL+L^2$. This shows Wasserstein embeddings are robust towards missing points in the input point cloud.

Table 3. Results of part segmentation on ShapeNetPart dataset with supervised (Sup) and self-supervised (S-Sup) pre-training. Bold represents the best results, and underline represents the second best.

Method	Pre-Train	mIoU	aero	bag	cup	car	chair	car phone	guitar	knife	lamp	laptop	motor	mug	pistol	rocket	skate board	table
CL+ L^2	Sup	77.61	73.54	62.76	72.41	64.75	82.42	65.41	88.95	<u>83.01</u>	75.83	93.52	43.03	83.25	74.87	46.01	62.39	77.54
	S-Sup	78.94	75.40	62.74	72.67	67.73	84.73	68.02	89.19	83.35	76.98	94.48	43.15	84.19	75.11	49.60	67.81	77.91
CL+Cosine	Sup	74.75	67.80	62.13	78.66	66.26	80.00	61.14	86.47	79.34	74.25	92.24	48.70	84.91	70.82	45.63	63.97	73.12
	S-Sup	78.49	72.51	<u>68.09</u>	71.44	68.35	84.47	63.38	88.69	80.30	75.94	94.45	48.81	<u>88.56</u>	74.08	47.37	<u>69.12</u>	77.90
WPCE	S-Sup	79.92	77.47	69.06	74.22	66.59	86.79	66.23	89.30	81.77	75.97	94.60	42.29	88.71	74.33	41.05	67.39	79.28
SSW-AE	S-Sup	75.20	68.83	59.61	69.65	64.23	81.98	62.00	86.92	80.27	73.70	92.82	38.46	85.12	68.26	43.97	60.65	73.69
CL+ SW_2 (Ours)	Sup	81.40	80.50	64.69	74.41	70.97	87.34	<u>69.71</u>	<u>89.34</u>	82.96	<u>77.59</u>	95.31	57.28	88.03	77.14	<u>53.18</u>	69.60	<u>79.84</u>
	S-Sup	<u>81.17</u>	<u>78.98</u>	66.90	<u>77.98</u>	<u>70.35</u>	<u>86.91</u>	70.57	89.39	82.85	77.99	<u>94.77</u>	<u>56.20</u>	87.18	<u>76.10</u>	53.98	69.60	80.10

4.5. 3D Shape Interpolation

We further examine the quality of our learnt space by performing shape interpolation between inter and intra class point cloud instances. The main aim of conducting this task is to examine which learnt space is capable of capturing geometric and structural information needed to generate consistent interpolations of 3D point clouds. One can inspect the quality of latent space by examining the smoothness of an interpolation path or the quality of the in-between generated samples in terms of noise. As interpolation is a synthesis task, we need a decoder network to reconstruct the object given its embedding. For this, we train an encoder-decoder network with our contrastive loss (Eq. 6) on the embeddings for the encoder, along with a reconstruction loss for the decoder. We use the encoder and decoder proposed by FoldingNet [24], which learns to deform a unit sphere and take the shape of a 3D object’s surface. We found that optimizing the network for learning discriminative features as well as detailed reconstruction is difficult. As our contrastive loss aims to pull point clouds closer with similar global representations, it becomes difficult to accurately reconstruct the input point cloud without fine-grain characteristic information. A simple way to deal with this issue is to assign weightage to the individual loss terms, with the weights summing to 1. In order to train an encoder-decoder network, the total effective loss is defined by taking a weighted sum of our contrastive loss and a reconstruction loss, with weights being 0.2 and 0.8, respectively. We use Chamfer distance as the reconstruction loss. Interpolation results are shown in Figure 4. We can see that the interpolations done using Wasserstein embeddings follow a smooth path with relatively less noisy points. For example, in Figure 4 (a), we can see that for Euclidean (in step 2 to 3), the interpolated sample suddenly takes the shape of a lamp. Whereas the interpolation path is smoother for Wasserstein. Similar trend can also be observed in Figure 4 (b), where for Euclidean, the source chair suddenly transforms (in step 2 to 3) to take the shape of the target chair, whereas in Wasserstein, legs of the chair smoothly morph to become the base of target chair. We evaluate the quality of interpolated samples

based on the noise present in them. For each interpolated sample, we compute the following noise measure

$$\mathcal{L}_{noise} = \sqrt{\frac{1}{N_m} \sum_{i=1}^{N_m} \left\| p_i - \frac{1}{n_e} \sum \mathcal{N}(p_i) \right\|^2} \quad (8)$$

where, $\mathcal{N}(p_i)$ gives a set of neighboring points for p_i that has a cardinality of $n_e = 20$. For a given point cloud, the noise measure computes neighborhood variation vectors around every point and aggregates their squared norm to give an overall score of smoothness. A lower noise measure implies that the point cloud has a smooth surface and has less noise. In Table 4, we report the noise measure of the interpolated samples shown in Figure 4. Most of the interpolated samples from Wasserstein space have lower values of noise measure. We provide more results for interpolation in our supplementary material.

Table 4. Noise measure for interpolation results shown in Figure 4. Bold values represent smoother surfaces having less noise. The values are scaled by a factor of 10^3 .

Figure	Method	Step 1	Step 2	Step 3	Step 4
4 (a)	Euclidean	30.2	35.6	36.7	35.4
	Wasserstein	29.6	32.1	32.6	31.4
4 (b)	Euclidean	29.3	29.9	30.2	28.0
	Wasserstein	29.7	29.1	29.4	27.7

4.6. Explainability

We investigate what makes Wasserstein embeddings perform better, as shown in the downstream tasks. We visualize and compare the features captured by Wasserstein embeddings and Euclidean embeddings in Figure 1. These features are called critical points, as shown by [15]. The embedding of a point cloud is completely determined by these subset of points. The embedding for a point cloud would be the same, as long as the set of critical points is

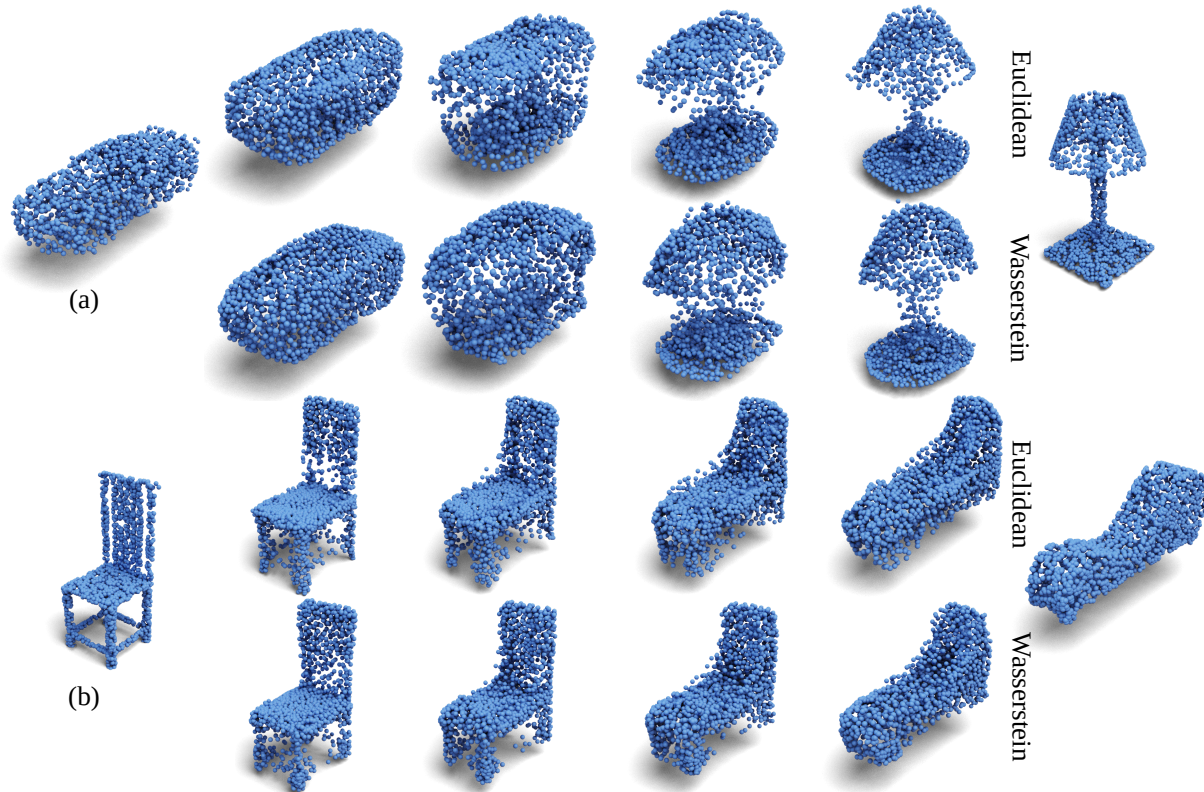


Figure 4. **Linear Interpolation** between source and target for two examples (a) Car to Lamp, (b) Chair to Chair from ShapeNet. The left and the right most represent original point clouds. The top row of each example shows results of reconstruction after interpolating two point cloud embeddings in Euclidean space. The bottom row of each example provides interpolation results in Wasserstein space. All rows follow the ratio of 0.8, 0.6, 0.4, and 0.2 (from left to right) with respect to the source.

unchanged. For a given point cloud, the critical points are those 3D points that contribute to the global embeddings after the max pooling layer. This implies that the number of critical points cannot be greater than that of the embedding size. The selection of critical points is extremely important, as they solely decide the embedding of a point cloud. This makes it clear that for good quality embeddings, critical points should best describe the given point cloud. In Figure 1, we can see that the network intelligently tries to summarize the point cloud by choosing boundary points as the critical points. Our Wasserstein embeddings are able to capture the full skeleton structure of the given point cloud, whereas critical points captured by Euclidean embeddings are comparatively poor with uneven distribution and missing parts. Thus, we can say that Wasserstein spaces are indeed better at preserving and capturing geometric structures amenable to the optimization task. We provide more examples in our supplementary material.

5. Conclusion

In this paper, we proposed to represent point clouds as discrete probability distributions in the Wasserstein space.

We built a contrastive learning method to learn Wasserstein embeddings for 3D point clouds. Our proposed method can be used as a pre-trained model in supervised and self-supervised settings for any downstream task. Empirically, we found that representations learnt using our pre-training of contrastive learning with Sliced Wasserstein distance captured the structure and underlying geometry better than standard Euclidean embeddings. With improved embeddings, our method outperformed all the existing methods, including our baseline with L^2 norm and Cosine similarity for all the downstream tasks (classification, segmentation, transfer learning, interpolation). We also show an interesting study of our self-explainable method by capturing critical points of point clouds better than embeddings in Euclidean space. For future work, a possible direction is to explore other related problems, such as domain adaptation for point clouds using optimal transport. Another interesting aspect is to consider complex datasets including multiple objects and scenes of point clouds.

References

- [1] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for contrastive learning of visual representations. In *International conference on machine learning (ICML)*, pages 1597–1607. PMLR, 2020. 3
- [2] S. Chopra, R. Hadsell, and Y. LeCun. Learning a similarity metric discriminatively, with application to face verification. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 539–546, 2005. 3
- [3] Nicolas Courty, Rémi Flamary, and Mélanie Ducoffe. Learning wasserstein embeddings. In *International Conference on Learning Representations (ICLR)*, pages 1–13, 2018. 2
- [4] Charlie Frogner, Farzaneh Mirzazadeh, and Justin Solomon. Learning entropic wasserstein embeddings. In *International Conference on Learning Representations (ICLR)*, 2019. 2
- [5] Meng-Hao Guo, Junxiong Cai, Zheng-Ning Liu, Tai-Jiang Mu, Ralph R. Martin, and Shi-Min Hu. PCT: point cloud transformer. *CoRR*, abs/2012.09688, 2020. 2
- [6] Keisuke Kawano, Satoshi Koide, and Takuro Kutsuna. Learning wasserstein isometric embedding for point clouds. In *International Conference on 3D Vision (3DV)*, pages 473–482, 2020. 2, 5
- [7] Prannay Khosla, Piotr Teterwak, Chen Wang, Aaron Sarna, Yonglong Tian, Phillip Isola, Aaron Maschinot, Ce Liu, and Dilip Krishnan. Supervised contrastive learning. *Advances in Neural Information Processing Systems (NeurIPS)*, pages 18661–18673, 2020. 3
- [8] Reza Mahjourian, Martin Wicke, and Anelia Angelova. Unsupervised learning of depth and ego-motion from monocular video using 3D geometric constraints. In *IEEE conference on computer vision and pattern recognition (CVPR)*, pages 5667–5675, 2018. 1
- [9] Daniel Maturana and Sebastian Scherer. Voxnet: A 3d convolutional neural network for real-time object recognition. In *IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 922–928, 2015. 1
- [10] Trung Nguyen, Quang-Hieu Pham, Tam Le, Tung Pham, Nhat Ho, and Binh-Son Hua. Point-set distances for learning representations of 3D point clouds. In *IEEE/CVF International Conference on Computer Vision (ICCV)*, 2021. 2, 5
- [11] Maximilian Nickel and Douwe Kiela. Poincaré embeddings for learning hierarchical representations. In *International Conference on Neural Information Processing Systems (NeurIPS)*, page 6341–6350, 2017. 2
- [12] Maximilian Nickel and Douwe Kiela. Learning continuous hierarchies in the lorentz model of hyperbolic geometry. In *International Conference on Machine Learning (ICML)*, pages 3776–3785. PMLR, 2018. 2
- [13] Linfei Pan, Ľubor Ladický, and Marc Pollefeys. Compression and completion of animated point clouds using topological properties of the manifold. In *International Conference on 3D Vision (3DV)*, pages 734–742, 2020. 1
- [14] Gabriel Peyré and Marco Cuturi. Computational optimal transport. *Foundations and Trends in Machine Learning*, pages 355–602, 2019. 3
- [15] Charles R. Qi, Hao Su, Kaichun Mo, and Leonidas J. Guibas. PointNet: Deep learning on point sets for 3D classification and segmentation. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017. 1, 7
- [16] Charles Ruizhongtai Qi, Hao Su, Matthias Nießner, Angela Dai, Mengyuan Yan, and Leonidas Guibas. Volumetric and multi-view CNNs for object classification on 3D data. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016. 1
- [17] Charles R. Qi, Li Yi, Hao Su, and Leonidas J. Guibas. PointNet++: Deep hierarchical feature learning on point sets in a metric space. In *International Conference on Neural Information Processing Systems (NeurIPS)*, page 5105–5114, 2017. 1, 4
- [18] Xie Saining, Gu Jiatao, Guo Demi, R. Qi Charles, Guibas Leonidas, and Litany Or. Pointcontrast: Unsupervised pre-training for 3d point cloud understanding. In *ECCV*, 2020. 2
- [19] Hang Su, Subhansu Maji, Evangelos Kalogerakis, and Erik Learned-Miller. Multi-view convolutional neural networks for 3D shape recognition. In *IEEE International Conference on Computer Vision (ICCV)*, pages 945–953, 2015. 1
- [20] Mikaela Angelina Uy, Quang-Hieu Pham, Binh-Son Hua, Duc Thanh Nguyen, and Sai-Kit Yeung. Revisiting point cloud classification: A new benchmark dataset and classification model on real-world data. In *International Conference on Computer Vision (ICCV)*, 2019. 5
- [21] Villani. *Topics in Optimal Transportation*. American Mathematical Society, 2003. 2
- [22] Yue Wang, Yongbin Sun, Ziwei Liu, Sanjay E. Sarma, Michael M. Bronstein, and Justin M. Solomon. Dynamic graph CNN for learning on point clouds. *ACM Transactions on Graphics (TOG)*, 2019. 1
- [23] Zhirong Wu, Shuran Song, Aditya Khosla, Fisher Yu, Linguang Zhang, Xiaoou Tang, and Jianxiong Xiao. 3D ShapeNets: A deep representation for volumetric shapes. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 1912–1920, 2015. 1, 4
- [24] Yaoqing Yang, Chen Feng, Yiru Shen, and Dong Tian. FoldingNet: Point cloud auto-encoder via deep grid deformation. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 206–215, 2018. 5, 7
- [25] Li Yi, Vladimir G. Kim, Duygu Ceylan, I-Chao Shen, Mengyan Yan, Hao Su, Cewu Lu, Qixing Huang, Alla Sheffer, and Leonidas Guibas. A scalable active framework for region annotation in 3D shape collections. *ACM Transactions on Graphics (TOG)*, 2016. 1, 5
- [26] Renrui Zhang, Ziyu Guo, Peng Gao, Rongyao Fang, Bin Zhao, Dong Wang, Yu Qiao, and Hongsheng Li. Point-m2ae: Multi-scale masked autoencoders for hierarchical point cloud pre-training. *arXiv preprint arXiv:2205.14401*, 2022. 2
- [27] Hengshuang Zhao, Li Jiang, Jiaya Jia, Philip H. S. Torr, and Vladlen Koltun. Point transformer. *CoRR*, abs/2012.09164, 2020. 2
- [28] Yuke Zhu, Roozbeh Mottaghi, Eric Kolve, Joseph J Lim, Abhinav Gupta, Li Fei-Fei, and Ali Farhadi. Target-driven visual navigation in indoor scenes using deep reinforcement learning. In *IEEE international conference on robotics and automation (ICRA)*, pages 3357–3364, 2017. 1