# PLURAL: 3D Point Cloud Transfer Learning via Contrastive Learning With Augmentations

Michael Biehler, Yiqi Sun, Shriyanshu Kode, Jing Li, and Jianjun Shi

*Abstract*— Unlocking the power of 3D point cloud machine learning models can be a challenge due to the need for extensive labeled datasets, which presents a challenge when applying these models to new domains. Transfer learning can help overcome this challenge by utilizing data from related tasks to enhance model performance. However, traditional (2D) transfer learning methods struggle with 3D point cloud domain adaptation, due to differences in physical environments and sensor configurations. To address this issue, we propose PLURAL, a novel 3D point cloud transfer learning methodology based on contrastive learning with augmentations. Our approach is inspired by the notion that high-level shape features are more transferable than low-level geometry features. We propose a co-training architecture that includes separate 3D point cloud models with domain-specific parameters, as well as a module for learning domain-invariant features. Additionally, PLURAL extends the approach of contrastive instance alignment to 3D point cloud modeling by considering physics-informed hard sample mining. Our experiments on simulation and real-world datasets demonstrate that PLURAL outperforms state-of-the-art transfer learning methods by a significant margin, effectively reducing the domain gap.

*Note to Practitioners*—The usage of 3D point cloud machine learning models is currently limited by the need for extensive labeled data. With our proposed framework, data from related tasks can be utilized to enhance the model performance on new applications or domains. PLURAL explicitly considers the acquisition of 3D point clouds by diverse sensors and in diverse environments. The method is highly adaptable and includes separate models with domain-specific parameters, making it applicable to a wide range of applications and domains.

*Index Terms*— Transfer learning, 3D point cloud, contrastive learning, hard sample mining.

## I. INTRODUCTION

**3D** POINT clouds have become increasingly important in real-world applications, such as predicting manufacturing quality and identifying objects in autonomous driving. The advancements in high-precision laser and LiDAR
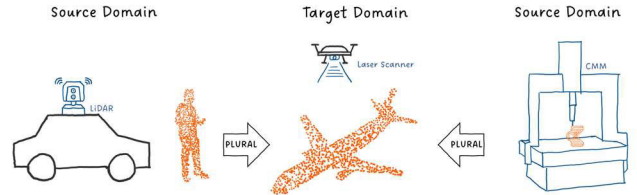
Fig. 1. Motivation of 3D Point Cloud Transfer Learning (PLURAL) approach.

sensors have revolutionized 3D point cloud modeling. However, obtaining labeled datasets for new domains remains a challenge, leaving many engineering problems unresolved. Different domains of 3D point cloud data are usually not independently and identically distributed, leading to domain shifts due to differences in physical environments or sensor configurations (e.g., different types of sensors, varying numbers of laser beams, and installation positions). One example of domain shift can occur due to differences in sensor configurations used to capture the same object and scene. As a result, most existing transfer learning methods are not readily applicable to this domain. To address these issues, we propose a transfer learning methodology that utilizes labeled data from related engineering problems to enhance the performance of an (unlabeled) target learning task based on contrastive learning with augmentations. Our method aims to effectively adapt 3D point cloud models from labeled source domains to a novel, unlabeled target domain by learning transferable features. As illustrated in Fig. 1, the goal of the PLURAL framework is to utilize 3D point cloud data from diverse source domains to enhance the predictive performance on a novel target domain. The 3D point clouds from different domains may exhibit significant domain shifts due to different acquisition sensors and object sizes and characteristics.

Domain shifts in 2D images typically manifest as changes in image appearance, such as blur, illumination, and weather conditions. In contrast, domain shifts in 3D point clouds are primarily caused by variations in shape geometry, resulting from both external factors and internal sensor configurations. Unlike 2D images, which have a uniform distribution of pixels, 3D point clouds can exhibit significant differences in low-level local geometry, making it challenging to transfer knowledge between different domains. Figure 2 serves as an illustration of this concept: in 2D scenes, such as transitions from the GTA [1] to Citiscapes [2] dataset, domain shifts are primarily characterized by changes in visual appearance. In contrast, 3D domain shifts predominantly manifest as variations in

**2D Domain Shift**

GTA Dataset [1]    Cityscapes Dataset [2]

**3D Domain Shift**

64-Beam LiDAR    32-Beam LiDAR

Same scene, same environmental conditions, **Different** sensor (KITTI Dataset [3])
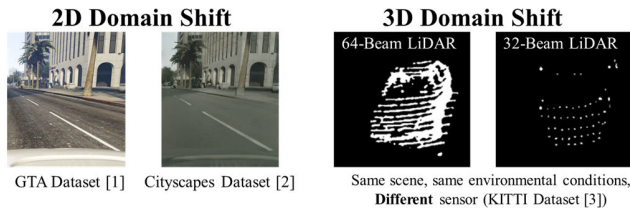
Fig. 2.    Illustration of characteristics of 2D and 3D domain shift.

geometry, stemming not only from external environmental factors but also from internal sensor configurations (e.g., differing laser beam configurations in datasets like KITTI [3]).

To address these challenges, we propose a novel transfer learning framework for 3D point cloud learning tasks called PLURAL. Our architecture includes separate models for each domain, along with a domain-agnostic learning module and a discriminative head. We hypothesize that high-level shape features are more transferable than low-level, local features, which form the basis of our architecture design. Our framework also incorporates a novel contrastive learning mechanism that leverages contrastive instance alignment, data augmentation, and hard sample mining to improve transferability and prevent the model from getting stuck in local minima. Unlike prior work in 3D point cloud transfer learning that used a self-training pipeline, PLURAL co-trains labeled data from multiple source domains and the target domain, leading to improved transfer of knowledge across domains.

The main contributions of this paper are as follows:

- We propose a novel transfer learning approach for 3D point cloud modeling that utilizes a unique architecture design involving parallel training of deep neural networks for a source and (unlabeled) target domain, with separate feature extractors.
- We introduce innovative techniques tailored for 3D point cloud data, including contrastive instance alignment and hard sample mining. For example, we employ cosine similarity distance, which possesses advantageous characteristics such as scale invariance for varying point densities and the preservation of directional information in spatial relationships, even when the scale of objects varies across domains.
- We conduct extensive simulations and case studies to demonstrate that our proposed method outperforms existing state-of-the-art transfer learning methods in terms of accuracy and generalizability.
- We provide new insights into how to encourage deep neural networks to learn transferable 3D point cloud features, with potential applications across various domains.

The structure of the article is as follows: In Section II, we provide a literature review on existing work in 3D point cloud transfer learning. Section III presents our proposed PLURAL framework for point cloud transfer learning via contrastive learning with augmentations. Section IV validates the effectiveness of the proposed methodology through extensive case studies and comparisons with existing benchmark methods. In Section V, we discuss the implications of our findings and outline potential avenues for future research. Finally, we conclude the article in Section VI.

## II. LITERATURE REVIEW

In this section, we provide a comprehensive review of the existing literature on transfer learning for 3D point cloud data. While there is a considerable amount of research on transfer learning with 2D images, the literature on transfer learning with 3D data is still limited.

### A. Methods Based on 2D Transfer Learning

Similar to other works in the field of 3D vision, a popular approach is to utilize well-established frameworks from 2D image processing and extend them to 3D data by either transforming the data itself (e.g., birds-eye-view) or extending network components to 3D (e.g., 3D convolutions). Along this direction, Imad et al. [4] convert the raw point cloud to a 2D bird's eye view and utilize a 2D convolutional neural network to perform semantic segmentation of 3D objects. Yan et al. [5] utilize deep learning models for existing sensor modalities (e.g. RGB images, 2D LiDAR) to learn a new 3D LiDAR-based human classifier from other sensors over time, taking advantage of a multisensor tracking system. Chai and Zhou [6] introduced a transfer learning methodology for industrial fault diagnosis, featuring a similarity learning-based discrimination module to identify fault prototypes (FPs) that are both representative of individual faults and discriminative across various fault categories. Additionally, a fault prototypical-adaptation module has been incorporated to adapt multiple FPs to the target dataset, enhancing category-specific domain invariance with precision. Chai et al. [7] introduced a multisource-refined transfer network to address fault diagnosis challenges in the presence of both domain and category inconsistencies. The network employs a multisource-domain-refined adversarial adaptation strategy to mitigate category-wise distribution inconsistencies within source–target domain pairs, avoiding negative transfer issues. Additionally, a multiple classifier complementation module leverages various diagnostic knowledge sources by transferring source classifiers to the target domain based on similarity scores, resulting in target-faults-discriminative and domain-refined-indistinguishable feature representations.

### B. Methods Based on 3D Transfer Learning

3D point cloud learning constitutes a highly active research field. Prior to the emergence of Point-Net [8], point cloud segmentation methods within the realm of deep learning typically relied on multi-view approaches [9] or volumetric techniques [10]. PointNet [8] marked a pivotal shift as the first deep learning-based method designed to learn directly from individual points. Utilizing point-wise multi-layer perceptrons, PointNet extracts global features, while its subsequent refinement, PointNet++ [11], was extended to encompass local information. PointConv [12] introduced point-wise convolution operators that convolute points with their neighboring counterparts. Edge-conditioned convolution

(ECC) treats each point as a graph vertex and employs graph convolution. RandLA-Net [13] adopts attention scores for points as soft masks, replacing the original pooling layer. Recently, transformer-based methods have gained attention, with point cloud transformer (PCT) [14] pioneering by incorporating self-attention layers into the original PointNet [8] framework.

However, annotating large volumes of real-world data for deep learning-based approaches in this domain is often challenging or even impractical. Recent methodologies have applied transfer learning directly within the 3D data domain. One viable solution involves the application of simulation-to-reality (sim2real) methods, where learning occurs with simulated data, and the acquired knowledge is subsequently transferred to real-world applications. This strategy, commonly employed in robot learning for tasks such as vision-related robotic tasks [15], [16], [17], typically involves rendering simulated scenes into RGB images, often with additional depth, thermal, or flow images. Deep learning-based neural networks are then pretrained using synthetic data and adapted to the real world through domain adaptation [18]. While this approach is prevalent in computer vision tasks involving images, limited work has been done regarding transfer learning on 3D point clouds [19]. Horache et al. [20] proposed utilizing a multi-scale U-Net-based method for descriptor matching for the specific learning task of 3D point cloud registration. Their approach allows the transfer of descriptors for registration on an unknown dataset without any supervision. Xiao et al. [21] focus on the transfer learning between synthetic and real (sim-to-real) 3D point clouds. Their approach translates synthetic point clouds to have a similar appearance and sparsity as real point clouds. However, their approach does not consider the performance of the translated point clouds on downstream tasks such as semantic segmentation. Wu et al. [19] extend one of the most widely used transfer learning techniques to 3D learning: Their 3D neural network model is firstly pre-trained on the synthetic data and then fine-tuned on the real-world data. Xie et al. [22] proposed a method for unsupervised pre-training on a large source set of a 3D scene to improve the performance on a small target dataset. Wei et al. [23] proposed a weakly supervised learning scheme for 3D point cloud scene semantic segmentation to reduce the labor and time cost for annotation on 3D datasets.

### C. Methods Based on 3D Hard Sample Mining

Hard sample mining is the process of selecting difficult training examples to improve model performance in challenging cases. It leads to better model generalization and performance by focusing the training on the most informative and hardest samples. For 3D point cloud data, Du et al. [24] proposed a self-contrastive learning framework with hard negative sampling based on nonlocal self-similarity, aiming at accurate point cloud representation learning in a self-supervised fashion. Yang et al. [25] proposed a self-supervised contrastive learning 3D classification model, that includes a confusion-prone classes mining module that mines classes with small inter-class variations. In the field of 3D vehicle detection, Zeng et al. [26] observed, that many region proposals contain no vehicles, or many region proposals contain simple examples, so online hard example mining [27] is adopted to augment the dataset.

To summarize, the development of 3D transfer learning is currently dependent on progress in 2D image transfer learning. However, to our knowledge, there has been no prior attempt to apply contrastive transfer learning directly to 3D architectures without relying on 2D data or network projections, while utilizing hard sample mining to examine areas of data that have been overlooked. This approach allows for the use of our transfer learning methodology in engineering applications that demand high precision, as using 2D projections would result in a loss of resolution and information.

### III. PLURAL METHODOLOGY

This section presents the PLURAL framework as an approach to point cloud transfer learning via contrastive learning with augmentations. We assume the following data scenario: From the target domain, a set of input point clouds $\mathcal{X}^T = \{X_i^T\}_{i=1}^{N^T}$ is available, where $i$ is the sample index, $N^T$ is the total number of samples, and the sample $X_i^T$ consists of a set of $n_i^T$ unstructured, varying-sized 3D measurement points (i.e., $X_i^T \in \mathbb{R}^{n_i^T \times 3}$). Note, that we do not assume that the response $\mathcal{Y}^T$ is available (i.e., unlabeled dataset). For the source domain, a set of input point clouds $\mathcal{X}^S = \{X_i^S\}_{i=1}^{N^S}$ is available along with a vector response $\mathcal{Y}^S = \{Y^S\}_{i=1}^{N^S}$, where $Y^S \in \mathbb{R}^{d_y^S}$, where $d_y^S$ denotes the dimension of the (multivariate) regression response or several one-hot labels of the classification response. Based on this dataset, we study the problem of unsupervised domain adaptation for 3D point cloud models, by adapting a 3D model $f_\theta$ parametrized by $\theta$ from a labeled source domain (i.e., $\{\mathcal{X}^S, \mathcal{Y}^S\}$) to an unlabeled target domain $\mathcal{X}^T$. The main objective of the PLURAL framework is to improve the performance of the 3D model $f_\theta$ on the unseen test set of the target domain, which requires careful consideration of the architecture and loss function to enable the learning of transferable features.

### A. Architecture Design and Big Picture

To achieve this goal, we employ 3D encoders for the target and source domains to learn domain-specific features. We then utilize a shared 3D encoder structure to learn transferable features, which are utilized in the discriminative head of the model. Fig. 3 provides a high-level overview of the proposed framework. Our hypothesis for the architecture design is that the 3D networks will gradually process domain-specific non-transferable features and acquire domain-invariant features. Therefore, we utilize domain-specific 3D encoders to learn low-level features that are distinct and specific to each data domain. Then we concatenate the features from the low-dimensional feature space and perform contrastive alignment for instance-level feature alignment. Finally, a discriminative head is utilized to obtain a supervision signal for the source domain and predict pseudo-labels for the target domain, which are subsequently updated during the joint optimization.
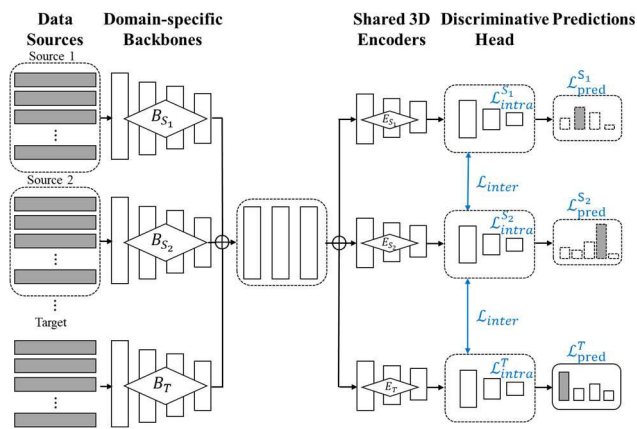
Fig. 3. Overview of the proposed PLURAL framework.

Transfer learning for 3D objects is challenging due to extreme geometry shifts, including density variations and different occlusion ratios of point clouds caused by diverse physical environments and sensor configurations. Unlike 2D domain models trained on ImageNet, 3D point cloud modeling lacks a transferable, well-trained backbone. One reason is the difficulty of reducing domain shifts in geometric representations for low-level features in the 3D model architecture. At a high level, we expect the 3D model to progressively process domain-specific non-transferable features and learn domain-invariant features. The neural network architecture facilitates this process by progressively learning domain-specific and domain-invariant information across its layers. This is a result of the network's hierarchical structure. In the early layers, particularly in convolutional neural networks (CNNs), convolutional layers utilize filters to capture local patterns, detecting simple features like edges and textures. Subsequent pooling layers then downsample spatial dimensions, highlighting salient features while reducing sensitivity to precise spatial locations. As the data advances through the network, deeper layers combine low-level features into more intricate and abstract representations. This gradual abstraction of features contributes to the network's capacity for invariance to irrelevant variations, fostering generalization. It is important to note that while conventional interpretations of invariance learning often treat domain-invariant features as separable from domain-specific features, our proposed approach emphasizes the network's learned ability to "filter out" domain-specific information. This process minimizes the influence of domain-specific details while retaining pertinent information for achieving domain invariance.

Our architecture design utilizes domain-specific 3D encoders that learn different mapping functions to convert unstructured 3D point clouds into a low-dimensional feature space. This leads to domain adaptation on the target domain while maintaining performance on the source domain, allowing for bi-directional knowledge sharing. Shared 3D encoders are then utilized to co-train with data samples from both domains, further compressing the outputs of the domain-specific 3D encoders. At this stage, we extend ideas from contrastive alignment learning in 2D vision, which encourages

the learning of domain-invariant features for deeper features that are more transferable because they have similar structures to grid-based feature maps in 2D image tasks. The discriminative head then predicts a regression response or classifies the 3D objects. Given labeled samples from the source domain, the detection head is trained to minimize a mean squared error (MSE) or cross-entropy loss for regression or classification tasks, respectively.

$$\mathcal{L}_{discr,reg} = \frac{1}{N^S} \sum_{i=1}^{N^S} \left\| Y_i^s - f_\theta\left(X_i^S\right) \right\|_2^2$$

$$\mathcal{L}_{discr,cls} = -\frac{1}{N^S} \sum_{i=1}^{N^S} \left\| Y_i^s \cdot \log\left(f_\theta\left(X_i^S\right)\right) \right\|_2^2 \quad (1)$$

### B. Contrastive Instance Alignment

In this section, we utilize the concept of contrastive instance alignment induced by pseudo-labels. The fundamental notion behind contrastive alignment is to minimize the feature distance between similar samples from different domains. We do this by leveraging pseudo-labels, which enable us to enhance the discriminative ability of the network and ensure that similar samples are aligned in the low-dimensional feature space. This approach significantly improves the generalization capability of the model and enhances its ability to tackle domain shift problems. Specifically, we choose the feature instance pair $(F_i^S, F_j^T)$ based on a similarity criterion as follows. Note that the similarity criterion slightly differs for regression and classification tasks. For each source feature instance $F_i^S$, we aim to find a feature instance $F_{j*}^T$ from the target domain that maximizes the cosine similarity:

$$j_{regr}^* = \max_{1 \le j \le N^T}\left\{\Phi(F_i^S, F_j^T)\right\}, \ 1 \le i \le N^S \quad (2)$$

$$j_{class}^* = \max_{1 \le j \le N_c^T}\left\{\Phi(F_i^S, F_j^T)\right\}, \ 1 \le i \le N_c^S, \ c = 1, .., |C|,$$

where $\Phi\left(F_i^S, F_j^T\right) = \frac{F_i^S \cdot F_j^T}{\left\| F_i^S \right\| \cdot \left\| F_j^T \right\|}$ calculates the cosine similarity between features of a source sample and a target candidate. For the classification task, $N_c^S$ and $N_c^T$ denote the total number of samples in class c in the source and target dataset (i.e., $S$ and $T$), respectively. $|C|$ denotes the total number of categories. In addition to minimizing the inter-class distance between domains, we also constrain the intra-class distance between different samples within the same domain. Hence, we get the following loss functions for the contrastive alignment depending on the discriminative task (regression versus classification).

$$\mathcal{L}_{inter, \ regr}(S, T) =$$

$$-\sum_{i \in N^S} = log \frac{\exp\left(F_i^S \cdot F_{j*}^T/\tau\right)}{\exp\left(F_i^S \cdot F_{j*}^T/\tau\right) + \sum_{j \in N^T} \exp\left(F_i^S \cdot F_j^T\right)}$$

$$\mathcal{L}_{inter, \ class}(S, T) =$$

$$-\sum_{c=1}^{|C|} \sum_{i \in N_c^S} log \frac{\exp\left(F_i^S \cdot F_{j*}^T/\tau\right)}{\exp\left(F_i^S \cdot F_{j*}^T/\tau\right) + \sum_{j \in N_{|C| \backslash c}^T} \exp\left(F_i^S \cdot F_j^T\right)}$$

$$(3)$$

$\mathcal{L}_{intra, \ regr}(D)$

$$= -\sum_{i \in N^D, j \in N^D} log \frac{\exp\left(F_i^D \cdot F_{j*}^D / \tau\right)}{\exp\left(F_i^D \cdot F_{j*}^D / \tau\right) + \sum_{j \in N^D} \exp\left(F_i^D \cdot F_j^D\right)},$$

$$D = \{S, T\}$$

$\mathcal{L}_{intra, \ class}(D)$

$$= -\sum_{c=1}^{|C|} \sum_{i \in N^D, j \ \in N^D} log \frac{\exp\left(F_i^D \cdot F_{j*}^D / \tau\right)}{\exp\left(F_i^D \cdot F_{j*}^D / \tau\right) + \sum_{j \in N_{|C|\backslash c}^D} \exp\left(F_i^D \cdot F_j^D\right)},$$

$$D = \{S, T\} \tag{4}$$

$$\mathcal{L}_{contr,align} = \mathcal{L}_{inter}(S, T) + \mathcal{L}_{intra}(S) + \mathcal{L}_{intra}(T), \tag{5}$$

where $\tau$ denotes a tuning parameter for the strength of domain adaption. The contrastive alignment loss $\mathcal{L}_{contr,align}$ considers the pairwise relations of samples between and within domains (source and target) to enable inter-domain transfer learning and improved discriminative performance on intra-domain tasks. Finally, by combining the loss terms in Equations 1 and 5, we optimize the model $f_\theta(\cdot)$ by minimizing the following loss:

$$\min_\theta \mathcal{L}_{discr} + \lambda \cdot \mathcal{L}_{contr,align} \tag{6}$$

where $\lambda$ is a tuning parameter to balance domain adaptation (i.e., contrastive alignment) and the learning of the discriminative task. However, since the features of point clouds are sparsely distributed, it is difficult to achieve effective alignment between domains by using global distribution alignment. In our experiments, we found that the simple use of contrastive alignment introduces the mismatch in point density and occlusion ratio between sample distributions of pseudo-labels and ground truths in the target domain. Therefore, we utilize effective augmentation via hard sample mining to further enhance domain adaptation.

### C. Hard Sample Mining

A straightforward use of the contrastive instance alignment tends to introduce the mismatch between the sample distributions obtained by pseudo-labels and ground truths on the target domain. First, pseudo-labels are more concentrated in the patterns with dense point clouds than those with sparse point clouds. Second, pseudo-labels cannot completely cover the patterns of severe occlusions. Therefore, most instances induced by pseudo-labels can be viewed as "easy samples" with sufficient points or complete geometry. However, we believe that the neglected "hard samples", which are more likely to be distributed in the tails of the predictive distribution, are equally important to 3D transfer learning.

Hard sample mining transforms point clouds by considering the specific geometry mismatches. There are some recent works [27], [28], concurrent to our own that select hard samples for training deep neural networks. Similar to our approach, all these methods based their selection on the current loss for each data point. More generally, the alternating steps that define a hard sample mining algorithm are as follows:

a) For some period of time, a *fixed model* is used to find new examples to add to the active training set,

b) Then, for some period of time, the model is trained on the *fixed* active *training set*.

To maximize the network learning, augmented sample $X_{i,hsm}^T$ generated by the augmentor should satisfy two requirements:

i. $X_{i,hsm}^T$ should be more challenging than $X_i^T$, i.e., we aim for the discriminative loss to be larger ($\mathcal{L}_{discr}(X_{i,hsm}^T) \geq \mathcal{L}_{discr}(X_i^T)$);

ii. $X_{i,hsm}^T$ should not lose its 3D shape features, meaning that it should describe a shape that is not too far (or different) from $X_i^T$.

To achieve requirement (i), a simple way to formulate the loss function for the hard sample mining (denoted by $\mathcal{L}_{hsm}$) is to maximize the difference between the discriminative losses on $X_i^T$ and $X_{i,hsm}^T$, or equivalently, to minimize

$$\mathcal{L}_{aug} = -\left(\mathcal{L}_{discr}(X_{i,hsm}^T) - \mathcal{L}_{discr}(X_i^T)\right). \tag{7}$$

Note that, for $X_{i,hsm}^T$ to be more challenging than $X_i^T$, we assume that $\mathcal{L}_{discr}(X_{i,hsm}^T) \geq \mathcal{L}_{discr}(X_i^T)$) and a larger $\mathcal{L}_{discr}(X_{i,hsm}^T)$ indicates a larger magnitude of augmentation, which can be defined as $\varepsilon = \mathcal{L}_{discr}(X_{i,hsm}^T) - \mathcal{L}_{discr}(X_i^T)$. However, if we naively minimize Eq. 7 for $\mathcal{L}_{aug} \rightarrow 0$, we encourage $\mathcal{L}_{discr}(X_{i,hsm}^T) - \mathcal{L}_{discr}(X_i^T) \rightarrow \infty$. This is due to the fact that $\mathcal{L}_{aug}$ minimizes the negative difference between the discriminative loss of the hard and original sample respectively (i.e., $X_{i,hsm}^T$ and $X_i^T$). This effectively maximizes the difference. Since the $\mathcal{L}_{discr}$ is fixed, without further constraints, we could choose the hard sample $X_{i,hsm}^T$ arbitrary far away from the original sample and improve our objective value. Hence, a naive solution for $X_{i,hsm}^T$ is an arbitrary, random shape regardless of $X_i^T$. Such a $X_{i,hsm}^T$ clearly violates requirement (ii). Hence, we further restrict the augmentation magnitude $\varepsilon$. Therefore, we upper-bound $\mathcal{L}_{discr}(X_{i,hsm}^T)$ with a dynamic parameter $\delta$:

$$\varepsilon = \mathcal{L}_{discr}(X_{i,hsm}^T) - \mathcal{L}_{discr}(X_i^T) \leq \delta \cdot \mathcal{L}_{discr}(X_i^T) = \varepsilon_{max}, \tag{8}$$

where $\varepsilon_{max}$ is the upper bound on the augmentation magnitude $\varepsilon$. Note that, when we learn the hard samples, the 3D encoders are fixed so $\mathcal{L}_{discr}(X_i^T)$ is fixed. Hence, $\varepsilon$ depends only on $\delta$. Since it should be non-negative, we thus ensure $\delta \geq 1$. Moreover, considering that the 3D models are very fragile at the beginning of the training, we pay more attention to training the classifier rather than generating a challenging hard sample $X_{i,hsm}^T$. Hence, $\varepsilon$ should not be too large at the beginning of the training process, meaning that $X_{i,hsm}^T$ should not be too challenging. Later, when the 3D model's discriminative ability becomes more powerful, we can gradually enlarge $\varepsilon$ to allow the augmentor to generate a more challenging $X_{i,hsm}^T$. Therefore, we design a dynamic $\delta$ to control $\delta$ with the following formulation:

$$\delta = \max\left(1, \frac{1}{\mathcal{L}_{discr}}\right), \tag{9}$$

where $\max(1, \cdot)$ ensures $\delta \geq 1$. At the beginning of the network training, the discriminate loss will be larger, since

the predictions may not be accurate. Hence, the inverse discriminative loss $\frac{1}{\mathcal{L}_{discr}}$ is generally small, resulting in a small $\delta$, and $\varepsilon$ will also be small according to Equation 8. When the classifier becomes more powerful, the discriminative loss $\mathcal{L}_{discr}$ will decrease, its inverse $\frac{1}{\mathcal{L}_{discr}}$ will increase, and we will have larger $\delta$ and $\varepsilon$ accordingly. Hence, the final hard sample loss is computed as follows,

$$\mathcal{L}_{hsm} = -\left(\mathcal{L}_{discr}\left(X_{i,hsm}^T\right) - \delta \cdot \mathcal{L}_{discr}\left(X_i^T\right)\right). \quad (10)$$

We here propose a novel algorithm to efficiently obtain hard samples, that optimizes the hard sample loss $\mathcal{L}_{hsm}$ without directly conducting gradient-based optimization to find the hard samples $X_{i,hsm}^T$. This allows us to find hard samples more efficiently. This method combines two major components: Firstly, it simulates object occlusions by altering the complete geometry of easy samples. Concretely, we calculate the viewpoint of a certain sample, randomly select a part of the viewpoint, and discard critical points on these angles. Secondly, it discards critical points (along the gradient direction) from existing dense point clouds. The exclusion of certain critical points is intentional and serves a specific purpose: We aim to train a model that learns more robust features, which can adapt to the challenges posed by 3D domain shifts, such as occluded objects or sparse measurement points. By systematically removing certain critical points, we encourage the model to focus on learning features that are resilient to geometric variations and domain shifts, ultimately promoting a more robust and adaptable representation. This also has a physical interpretation: this process simulates the change in the number of laser beams among different 3D sensors. The idea is to use the information of the gradient of the neural network model which tells us for each 3D input point whether the model performance will increase if we take a small step in the gradient direction.

The core component of hard sample mining is the attribution score: The attribution value assigns each point a score reflecting its contribution to the discriminative model loss. Aggregations of highly scored points indicate important segments/subsets in a 3D point cloud. If a measurement point with high attribution scores is discarded, the model performance decreases significantly. Therefore, discarding a relatively large number of points with high attribution scores leads to a new sample that is a "hard sample" for the model to predict. Fig. 4 illustrates the hard sample mining algorithm, which includes the random selection of a viewpoint, the calculation of point attributions of the points in this viewpoint, and the subsequent deletion of points with large attribution scores until the termination criterion $\mathcal{L}_{discr}\left(X_{i,hsm}^T\right) - \mathcal{L}_{discr}\left(X_i^T\right) > \delta \cdot \mathcal{L}_{discr}\left(X_i^T\right)$ is met.

The transformed point clouds focus on effective contrastive instance alignment by reducing the distribution mismatch of the target domain induced by pseudo-labels.

Figure 5 depicts the PLURAL procedure for constructive alignment. Initially (Fig. 5a), there exists a feature space mismatch between the source and target domains due to domain shifts, including differences in sensor configurations. This mismatch results in inaccurate predictions (red points) in areas of the target domain that lack a strong match to the
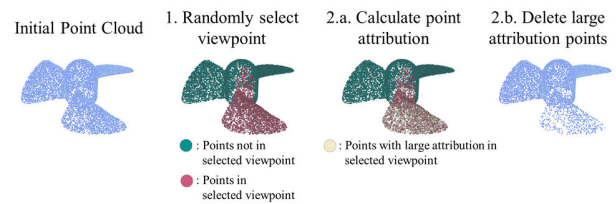


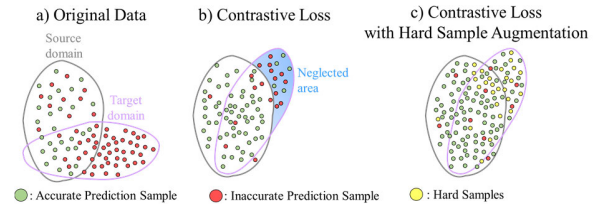Fig. 4.    Illustration of hard sample mining algorithm.



Fig. 5.    PLURAL procedure for contrastive alignment with hard sample augmentation.

source domain. While the application of the contrastive loss improves domain alignment to some extent, it tends to align easily recognizable 3D objects but often overlooks challenging samples (blue shaded area, Fig. 5b) characterized by severe occlusions or density variations. Consequently, the transfer learning model may continue to produce inaccurate predictions in these regions, leading to discrepancies in point density and occlusion ratios between the pseudo-labels and ground truths in the target domain. To address this challenge, we have developed a dedicated hard sample mining algorithm (Fig. 5c) that transforms point clouds while taking into account specific geometry mismatches, such as severe occlusions or density variations. This approach results in a substantial enhancement in domain alignment, ultimately improving the overall effectiveness of our transfer learning framework.

The purpose of the hard sample mining is to further enhance the contrastive alignment scheme of the PLURAL framework. The original contrastive alignment focuses more on the alignment of easy 3D objects rather than the easily neglected hard samples with severe occlusions or density variations. To assess the impact of hard sample mining on alignment, we conducted an ablation study using a partial model variant, denoted as "PLURAL w/o HSM," which does not include hard sample mining. The findings from this analysis will be discussed in the experimental section.

### D. Unified PLURAL Framework

We propose a stepwise training procedure with a warm-up process to train the PLURAL framework as shown in Algorithm 2. Specifically, we first pre-train a source model on the labeled source domain and use it to generate pseudo-labels on the target set. We then conduct hard sample mining (Algorithm 1) and augment the target set. Next, we warm up the 3D model following Equation 6, which allows a more stable convergence in the early stages of training. For the remaining epochs, we update the pseudo-labels using stepwise co-training. During this process $f_\theta(\cdot)$ gradually adapts to the target domain while maintaining the in-domain performance.

---

**Algorithm 1** Hard Sample Mining Algorithm

**Inputs:**
- "Easy" Target sample: $X_i^T$
- Dynamic upper-bounding parameter $\delta$
- Current 3D model $f_\theta$ and associate discriminative loss $\mathcal{L}_{discr}$

**Initialize:** Number of discarded points $n_{disc} = 0$
**Output:**
- "Hard" Target sample: $X_{i,hsm}^T$

**Algorithm:**
1. Select a random viewpoint $X_i^{T,view}$ of the 3D point cloud $X_i^T$
2. Discard points along the gradient direction in the selected viewpoint

while $\mathcal{L}_{discr}\left(X_{i,hsm}^T\right) - \mathcal{L}_{discr}\left(X_i^T\right) \le \delta \cdot \mathcal{L}_{discr}\left(X_i^T\right)$

    a. Calculate the attribution $a_j$ of each point $j$ in the viewpoint of the 3D point cloud $X_{i,j}^{T,view}$ via $a_j = X_{i,j}^{T,view} \nabla M_\theta(X_{i,j}^T)$, where $j = 1, \dots, n_p$ and $n_p$ is the number of measurement points in the 3D point cloud.

    b. Successively discard point $j$ with the largest attribution $a_j$, $n_{disc} = n_{disc} + 1$

if $n_{disc} \ge \frac{1}{3}\left|\left\{X_i^{T,view}\right\}\right|$ (do not want to discard too many points in viewpoint)

    c. Go back to Step 1, and select a new viewpoint

---

**Algorithm 2** PLURAL Algorithm

**Inputs:**
- **Source:** Labeled point cloud dataset from source domain $D^S = \left\{(X_i^S, Y_i^S)\right\}_{i=1}^{N_S}$
- **Target:** Unlabeled input point cloud from target domain $\left\{X_i^T\right\}_{i=1}^{N_T}$
- **Algorithm parameters:** Network architecture and termination tolerance $\epsilon$

**Output:**
- Learned network weights $\theta$ of model $f(\cdot)$

**Algorithm:**
1. Pretrain 3D base model
   $f^{init} = \text{fit}(D^S)$
2. Generate pseudo-labels for target domain samples
   $\left\{\overline{Y_i^T}\right\}_{i=1}^{N_T} = \text{predict}\left(f^{init}, \left\{X_i^T\right\}_{i=1}^{N_T}\right)$
3. Mine hard samples to augment the target set
   $D_{hsm_0}^T = \left\{\left(X_{i,hsm_0}^T, \overline{Y_{i,hsm_0}^T}\right)\right\}_{i=1}^{N_{T,hsm_0}} = \text{hsm}\left(\left\{\left(X_i^T, \overline{Y_i^T}\right)\right\}_{i=1}^{N_T}\right)$
4. Initialize the model with 3D base model
   $f_\theta = f^{init}$
5. Warm start of PLURAL model
   $f_\theta^0 = \text{fit}(D^S, D_{hsm_0}^T)$
6. PLURAL iteration: Iteration index $k$
   While not converged:
   *6.1 Update pseudo labels:*
   $\left\{\overline{Y_{i,hsm_k}^T}\right\}_{i=1}^{N_{T,hsm_k}} = \text{predict}\left(f_\theta^0, \left\{X_i^T\right\}_{i=1}^{N_T}\right)$
   *6.2 Add new hard samples to the target dataset*
   $D_{hsm_k}^T = \left\{\left(X_{i,hsm_k}^T, \overline{Y_{i,hsm_k}^T}\right)\right\}_{i=1}^{N_{T,hsm_k}}$
   $= \text{hsm}\left(\left\{\left(X_{i,hsm_{k-1}}^T, \overline{Y_{i,hsm_{k-1}}^T}\right)\right\}_{i=1}^{N_{T,hsm_{k-1}}}\right)$
   *6.3 Model update*
   $f_\theta^k = \text{fit}(D^S, D_{aug,k}^T)$
   Termination check: $f_\theta^{k-1} - f_\theta^k \le \epsilon$

---

In our transfer learning approach, the generation of pseudo-labels serves as a crucial bridge for knowledge transfer between the source and target domains. These pseudo-labels facilitate effective domain alignment and are essential for model adaptation. We prioritize adaptability and robustness, enabling the model to excel in both source and target domains, particularly in scenarios with challenges like occlusions and sparse data points. To mitigate negative transfer, we employ two key strategies: contrastive instance alignment with cosine similarity and hard sample mining. Pseudo-labels play a pivotal role in these strategies, promoting the alignment of semantically similar instances while fostering generalization and robust feature learning. Our approach aims to maintain a balance between adaptability and domain-specific preservation, ultimately enhancing generalizability and target domain performance.

*E. Hyperparameter Tuning*

We note that the use of machine learning algorithms commonly involves careful tuning of learning parameters requiring expert experience, rules of thumb, or brute force search. On the contrary, we view this issue as the global derivative-free optimization of an unknown (nonconvex) black-box function and utilize the Bayesian optimization procedure proposed by Snoek et al. [29] with its accompanying Python package "Spearmint" to automatically optimize the performance of the PLURAL algorithm for a given problem. This process is fully automated and can be parallelized for computational efficiency during training time. Bayesian optimization has been shown to outperform other global optimization algorithms for tuning parameter selection on several multimodal black-box functions [30].

## IV. SIMULATION STUDIES

In this section, we evaluate the PLURAL approach with simulated unstructured 3D point clouds against three benchmark methods. The data characteristics are as follows: (i) the parts are measured and represented by unstructured point clouds, and (ii) the goal is to develop a model to link the point cloud inputs with a scalar regression response. In the simulation studies, we will use conic shapes to predict the roundness error. In manufacturing applications, conic shapes are commonly used as reference objects for problems such as part-to-part variation pattern identification [31] or process control [32]. Therefore, we simulate truncated cone point clouds with $n_p$ (i.e., $i = 1, \dots, n_p$) following this procedure:

1. Generate random angles $\theta_i \in [0, 2\pi]$ and radius on the bottom ($r_1 = N(3, 2)$) and top ($r_2 = N(10, 2)$) of the cone;
2. Generate random coordinates $z_i^{init} \in [0, h^{init}]$, where $h^{init}$ is the height of the cone;
3. Add random perturbations to the radius $r_i$ and height $h_i$:

$$r_i^\Delta = N(0, 0.5); \quad h_i^\Delta = N(0, 0.1 \cdot h^{init}); \quad h_i = h_{init} + h_i^\Delta; \quad z_i = z_i^{init} + h_i; \tag{11}$$

4. Add perturbations to the wall coordinates of the cone:
   Wave amplitude $w_i^A = N(0.5, 1)$; Wavelength $w_i^A = N(5, 1)$; Wave phase $w_i^A = U[0, 2\pi]$; Eccentricity $E_i = (r_1 - r_2)/2$, Eccentricity perturbation

   $$E_i^\Delta = N(0, 0.2)$$

Fig. 6. Examples of generated truncated cones.

$$r_i^p = \frac{r_1 + r_2}{2} + E_i \cdot \left(1 - \frac{z_i}{h_i}\right) + r_i^\Delta + E_i^\Delta$$
$$+ w_i^A \cdot \sin\left(2\pi \cdot \frac{z}{w_i^A} + w_i^A\right)$$
$$+ w_i^A \cdot \cos\left(2\pi \cdot \frac{z}{w_i^A} + w_i^A\right) \tag{12}$$

5. Calculate the corresponding $x$ and $y$ coordinates on the surface of the cylinder:

$$x_i = r_i^p \cdot \cos(\theta_i); \ y_i = r_i^p \cdot \sin(\theta_i) \tag{13}$$

6. Add random noise to each coordinate to simulate measurement errors

$$X_i = (x_i + \varepsilon_x, \ y_i + \varepsilon_y, z_i + \varepsilon_z), \tag{14}$$

where $\varepsilon$ is an additive Gaussian noise centered at the respective coordinate (i.e., $\varepsilon_x = N(x_i, \sigma^2)$, $\varepsilon_y = N(y_i, \sigma^2)$, $\varepsilon_z = N(z_i, \sigma^2)$). Here, $x$ and $y$ are the Cartesian coordinates on the surface of the cylinder, given by the radius and the angle at each point. The resulting point cloud $\mathcal{X} = \{X_i\}_{i=1}^{n_p}$ will consist of $n_p$ 3 D points with $(x, y, z)$ coordinates that lie on the surface of the cylinder, with a wave and eccentricities added to the walls of the cylinder. Two examples of the generated truncated cone point clouds with $n_p = 3000$ are shown in Fig. 6.

Domain shift can occur in 3D point clouds due to different sensor configurations used for data acquisition. This means that the statistical properties of the data may differ between two datasets acquired using different sensors, even if they represent the same scene or object. Such variations can lead to reduced performance of downstream tasks.

To simulate different sensor configurations, we generate two datasets by applying subsampling techniques:

- *Dataset A: Uniform acquisition: Ground LiDAR*

Large-scale point clouds acquired by ground laser (LiDAR) scanning typically exhibit a uniform sampling structure [33], [34]. Therefore, we subsample a generated point cloud $\mathcal{X}$ to $n_p = 1500$ using uniform subsampling to obtain a sample $\mathcal{X}^A$ of dataset A.

-*Dataset B: Space filling acquisition: High-resolution terrestrial laser scanning (TLS)*

High-resolution terrestrial laser scanning (TLS) produces point clouds with billions of measurement points posing significant computational challenges. Therefore, it is common to apply postprocessing to reduce the size of the point clouds but preserve their spatial information. Methods such as 3D Hilbert curves [35] or regular 2D grid sorting [36] lead to 3D point
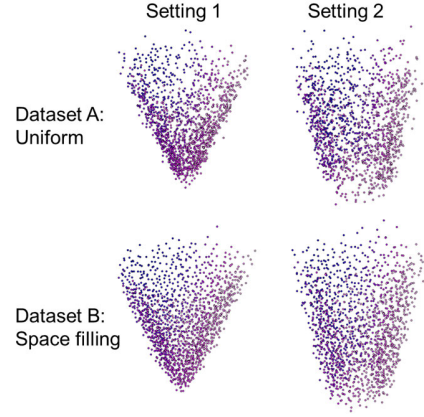


Fig. 7. Example of Datasets A and B mimicking different 3D sensor configurations.

clouds with space-filling properties. Therefore, we subsample a generated point cloud $\mathcal{X}$ to $n_p = 1500$ using furthest-point subsampling to obtain a sample $\mathcal{X}^B$ of dataset B.

Examples of those two different point cloud acquisition techniques with the same random generation settings (for comparison purposes) are visualized in Fig. 7. We can see that dataset A has a varying point density across the 3D surface, due to independent uniform sampling. On the other hand, dataset B has an even distribution on the 3D surface due to the consideration of space-filling properties for sampling.

An important quality characteristic of conic or cylindrical shape applications, such as hot steel rolling [37] or orthopaedical implants [38], is the roundness error of the manufactured parts. However, the application of the minimum zone tolerance (MZT) method, which provides the most accurate estimation of the roundness error, is computationally intensive [39]. However, the recently proposed ANTLER framework [40] can directly estimate such quality characteristics based on a point cloud input. Therefore, we will utilize the MZT method to compute the quality response in this simulation study via a nonlinear function in the following way.

Given a circumferential line $r(x, y, \phi)$, the roundness error $R(x, y)$ is defined by

$$R(x, y) = OC(x, y) - IC(x, y), \ (x, y) \in E_{r(x, \ y, \phi)}, \tag{15}$$

where $OC(x, y)$ and $IC(x, y)$ are the radii of the reference circles of center $(x, y)$ derived by $n$ finite measurement points and $E_{r(x, \ y, \phi)}$ is the area enclosed by $r(x, y, \phi)$:

$$OC(x, y) = \max_{\phi_i = i \times \frac{2\pi}{n}, \ i=1,\dots,n} r(x, y, \phi_i) \tag{16}$$

$$IC(x, y) = \min_{\phi_i = i \times \frac{2\pi}{n}, \ i=1,\dots,n} r(x, y, \phi_i) \tag{17}$$

Finally, the average roundness $(\overline{R}(x, y)$ over circumferential lines is used as the response $\mathcal{Y}$. Following this procedure, we generate $N = 100$ point cloud samples for each dataset. In particular, we obtain the source datasets $D_A^S = \left\{\mathcal{X}_j^A, \mathcal{Y}_j\right\}_{i=1}^N$ and $D_B^S = \left\{\mathcal{X}_j^B, \mathcal{Y}_j\right\}_{i=1}^N$. If the respective datasets serve as target datasets, the label is only utilized for validation purposes (i.e., $D_A^T = \left\{\mathcal{X}_j^A\right\}_{i=1}^N$ and $D_B^T = \left\{\mathcal{X}_j^B\right\}_{i=1}^N$).

## A. Benchmark Methods

For this simulation study, the authors of the ANTLER method [40] compared their proposed framework against six benchmarks and established state-of-the-art performance in the non-transfer learning setting on the conic shape dataset. Therefore, we use ANTLER as the basic, *non-transfer learning benchmark*. Furthermore, ANTLER serves as the domain-specific 3D encoder of our transfer learning approach. We note, that to the best of our knowledge, there exists no current method that can achieve transfer for 3D point cloud regression tasks to an unlabeled target domain. However, for comparison purposes, we still consider the following benchmarks, even though they have an "unfair" advantage due to their access to target labels:

- *Pretraining*

One of the most widely used transfer learning techniques is pretraining: Following Yan et al. [5], the ANTLER model is firstly pre-trained on the source dataset and then fine-tuned on the target dataset. The epochs to conduct fine-tuning are chosen based on binary search.

- *PointAugment*

Since PLURAL uses hard sample mining for augmentation, we compare our method against an advanced augmentation algorithm for 3D point cloud data, which also assumes access to target labels and is only trained on the target dataset. The goal of PointAugment [28] is to generate new point cloud samples by augmenting current samples to enrich data diversity. Instead of using fixed augmentation strategies, this model develops an augmentor that is trained alongside the model. Uses sample-aware data augmentation which regresses a specific augmentation function for each input sample by considering its geometric structure. The augmentor employs an adversarial learning strategy to optimize the augmentor so that the augmentor learns to produce samples that will best contribute to the model performance.

- *Unsupervised Domain Adaptation Methods*

To address transfer learning, another noteworthy category involves unsupervised domain adaptation methods. Consequently, we evaluate our approach against three state-of-the-art methods within this domain: CORAL (CORrelation ALignment) [41], LocIT (Localized Information Transfer) [42], and TCA (Transfer Component Analysis) [43]. These methods play a pivotal role in adapting models to new domains without the need for labeled data, emphasizing the importance of robust techniques in scenarios where source and target domains differ.

- *Multidomain Learning*

Different domains often contain information that can mutually benefit one another, a phenomenon frequently observed in ensemble learning. Multi-domain learning (MDL) presents a solution by harnessing domain information to augment the learning process. Therefore, we compare our method with two state-of-the-art methods in the field of point cloud MDL—Multi-Domain Knowledge Transfer (MDKT) [44] and Pointcloud Domain Adaptation Network (PointDAN) [45].

We stress, again, that we include *Pertaining* and *PointAugment* as non-transfer learning (TL) benchmarks, but those

TABLE I
SIMULATION STUDY PREDICTION RESULTS (BEST MODEL IN BOLD)

| Method/Dataset | $A \to B$ | | Method/Dataset | $B \to A$ | |
|---|---|---|---|---|---|
| $\delta = 0.1$ | RMSE | SD | $\delta = 0.1$ | RMSE | SD |
| ANTLER (no TL, on B only) | 0.183 | 0.037 | ANTLER (no TL, on A only) | 0.207 | 0.049 |
| Pretraining (with target label) | 0.179 | 0.043 | Pretraining (with target label) | 0.188 | 0.044 |
| PointAugment (no TL) | 0.162 | 0.028 | PointAugment (no TL) | 0.173 | 0.038 |
| CORAL [41] | 0.172 | 0.049 | CORAL [41] | 0.184 | 0.048 |
| LocIT [42] | 0.168 | 0.040 | LocIT [42] | 0.175 | 0.054 |
| TCA [43] | 0.178 | 0.046 | TCA [43] | 0.166 | 0.047 |
| PointDAN [45] | 0.154 | 0.048 | PointDAN [45] | 0.177 | 0.044 |
| MDKT [44] | 0.151 | 0.034 | MDKT [44] | 0.161 | 0.046 |
| PLURAL (w/o HSM) -ANTLER | 0.142 | 0.024 | PLURAL (w/o HSM) -ANTLER | 0.147 | 0.032 |
| **PLURAL-ANTLER** | **0.112** | **0.013** | **PLURAL-ANTLER** | **0.133** | **0.019** |
| $\delta = 1$ | RMSE | SD | $\delta = 1$ | RMSE | SD |
| ANTLER (no TL, on B only) | 0.212 | 0.050 | ANTLER (no TL, on A only) | 0.257 | 0.068 |
| Pretraining (with target label) | 0.208 | 0.049 | Pretraining (with target label) | 0.236 | 0.049 |
| PointAugment (no TL) | 0.194 | 0.035 | PointAugment (no TL) | 0.235 | 0.039 |
| CORAL [41] | 0.198 | 0.059 | CORAL [41] | 0.244 | 0.076 |
| LocIT [42] | 0.184 | 0.053 | LocIT [42] | 0.223 | 0.064 |
| TCA [43] | 0.193 | 0.056 | TCA [43] | 0.247 | 0.066 |
| PointDAN [45] | 0.186 | 0.048 | PointDAN [45] | 0.203 | 0.061 |
| MDKT [44] | 0.189 | 0.049 | MDKT [44] | 0.219 | 0.057 |
| PLURAL (w/o HSM) -ANTLER | 0.182 | 0.030 | PLURAL (w/o HSM) -ANTLER | 0.197 | 0.033 |
| **PLURAL-ANTLER** | **0.171** | **0.026** | **PLURAL-ANTLER** | **0.194** | **0.026** |

methods require access to the target labels, giving them an "unfair" advantage over our unsupervised domain adaptation framework.

## B. Prediction Results in the Simulation Study

We compare the proposed PLURAL method with the benchmark methods based on the Root Means Squared Error (RMSE) calculated at different levels of noise $\delta$. Table I reports the average and standard deviation of RMSE obtained via 10-fold cross-validation for simulation studies.

As shown in Table I, PLURAL outperforms all compared models by large margins. Especially on the transfer from the space-filling dataset B to the uniformly sampled dataset A, which exhibits a sparse shape representation PLURAL effectively closes the domain gap. Even though the pretraining and PointAugment benchmarks have access to the target labels, while PLURAL does not, they fail to achieve good model performance due to the small sample size and complex, high-dimensional data characteristics. In contrast, although starting with low-quality pseudo labels, PLURAL still outperforms those methods due to the effective co-training that incorporates labeled source data and augmented hard samples. In our study, we conducted comparisons with conventional unsupervised transfer learning methods, including CORAL, LocIT, and TCA. Our findings indicate that classical transfer learning methods face specific challenges when applied to 3D point clouds, owing to the distinctive characteristics of this data type. Notably, 3D point clouds exhibit an unstructured nature, often lacking predefined structures or topological awareness among points. Each point within an unstructured point cloud operates independently, and the distances to neighboring points vary,

This article has been accepted for inclusion in a future issue of this journal. Content is final as presented, with the exception of pagination.

10                                                                                              IEEE TRANSACTIONS ON AUTOMATION SCIENCE AND ENGINEERING

setting it apart from structured data. Additionally, the distribution of points can be irregular, with regions characterized by both sparsity and density, further complicating the application of traditional transfer learning approaches. We additionally compared our model with two state-of-the-art multi-domain transfer models, namely PointDAN [45] and MDKT [44], both of which exhibit better performance compared to unsupervised models. However, they still struggle with distribution mismatch issues. This becomes evident when contrasting them with our model without hard sample mining (PLURAL w/o HSM). Essentially, the alignment performance of the multi-domain models becomes comparable to the PLURAL w/o HSM, however, they may impose additional requirements on the datasets. However, our full PLURAL model, which additionally incorporates hard sample mining, proves effective in mitigating the remaining mismatch, resulting in superior prediction results.

The overall results validate the transferability of PLURAL on unsupervised domain adaptation benchmarks. Furthermore, we can see that with increasing noise levels the benefits of transfer learning decrease.

## V. CASE STUDY: SEWER DEFECT DETECTION

In this section, we evaluate PLURAL on a real-world dataset for transfer learning on 3D point clouds. In particular, we utilize the public AAU Sewer Defect Point Cloud dataset [46]. The dataset is focused on classifying defects in sewer pipes. The majority of the dataset consists of semi-synthetic data, while some point clouds of real sewer pipes were also recorded. The semi-synthetic dataset includes dry plastic pipes with defects introduced by displacing the pipes and placing rubber rings or bricks inside the pipes. Semi-synthetic in this setting means, that it is not obtained from real sewer pipes but still consists of point clouds collected from real 3D objects via time-of-flight sensors. Time-of-flight sensors are devices that emit light and measure reflection from an object to calculate distances to objects in the environment. These sensors are commonly used for 3D point cloud acquisition due to their ability to quickly capture depth information with high accuracy. For the target dataset with real pipes, a laser scanner is utilized leading to different data characteristics of the 3D point clouds with regard to the space-filling properties and sparsity patterns. Readers interested in the details of the dataset are referred to Haurum et al. [46]. Four different point clouds from both the source and target datasets are visualized in Fig. 8.

We can see that the properties of the point clouds differ across domains. The target domain consists of more discontinuous but denser shape surfaces. The source domain on the other hand consists of evenly distrusted but sparse 3D measurement points showing the significant difference between acquisition devices and environments across domains.

For the case study, the domain-specific 3D encoder is chosen as the PointNeXt model, which achieves state-of-the-art performance by utilizing advanced hyperparameter tuning for the popular PointNet model. Furthermore, this allows for a fair comparison with the non-transfer learning benchmark: The model used in the dataset paper [46] is also a PointNet model.
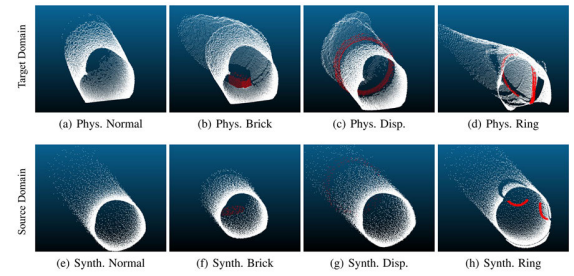


Fig. 8. Examples of point clouds from the source dataset (first row) and target dataset (second row) [46].

TABLE II
CASE STUDY PREDICTION RESULTS(BEST MODEL IN BOLD)

| Method/Metric | Accuracy | Precision | Recall | F1-Score |
|---|---|---|---|---|
| AAU Sewer [46] (with Pretraining) | 0.319 (0.062) | 0.321 (0.063) | 0.298 (0.138) | 0.229 (0.078) |
| Pretraining TL (with target label) | 0.415 (0.048) | 0.443 (0.047) | 0.421 (0.053) | 0.432 (0.049) |
| PointAugment (no TL) | 0.498 (0.040) | 0.138 (0.043) | 0.249 (0.046) | 0.167 (0.047) |
| CORAL [41] | 0.218 (0.031) | 0.224 (0.033) | 0.237 (0.031) | 0.229 (0.030) |
| LocIT [42] | 0.349 (0.042) | 0.227 (0.059) | 0.216 (0.032) | 0.218 (0.031) |
| TCA [43] | 0.328 (0.021) | 0.269 (0.024) | 0.266 (0.019) | 0.267 (0.021) |
| PointDAN [45] | 0.495 (0.069) | 0.139 (0.034) | 0.255 (0.019) | 0.177 (0.025) |
| MDKT [44] | 0.378 (0.017) | 0.286 (0.019) | 0.282 (0.017) | 0.284 (0.018) |
| PLURAL (w/o HSM) | 0.509 (0.042) | 0.533 (0.053) | 0.524 (0.045) | 0.547 (0.046) |
| **PLURAL** | **0.563 (0.032)** | **0.585 (0.041)** | **0.577 (0.031)** | **0.582 (0.037)** |

In total, there are 8100 source samples and 415 target samples. We perform 10-fold cross-validation to obtain the average classification performance on the target dataset. We evaluate the models by considering the accuracy, precision, recall, and F1-score and report the results in Table II.

From the results, it is evident that the PLURAL framework consistently outperforms the other benchmarks. We note that to the best of our knowledge, PLURAL is the currently best-performing model on the public AAU Sewer Dataset. We can see from the precision and recall scores, that the PLURAL framework has a well-rounded prediction performance and does not consistently misclassify a certain class more frequently. Classical unsupervised transfer learning methods, including CORAL, LocIT, and TCA, exhibit limited performance in the context of 3D point cloud data. Their challenges arise from their inability to adequately address the unique characteristics of this data type. A similar trend is evident in the performance of multi-domain models. Although there is an enhancement in accuracy, there is a deficiency in precision, recall, and F1 score. This deficiency indicates that these strategies struggle to achieve effective alignment and are more inclined to classify the more frequently occurring classes, an undesirable characteristic for a classifier. While these models demonstrate improvements compared to unsupervised models, they still inadequately address distribution mismatches within sparsely populated and occluded point patterns. For the PLURAL w/o HSM, we can see that the

contrastive alignment of the PLURAL framework is more effective than other benchmarks, but our full model with HSM can better exploit and match the information for both domains and achieves the best prediction performance. In this case study, 3D point clouds present significant variations in sparsity and point distribution between the source and target domain. Moreover, the presence of real-world debris and occlusions in the target domain poses obstacles to effective knowledge transfer, further impeding the performance of these methods. Even though the Pretraining and PointAugment methods have access to the target labels, they are not able to achieve a good performance due to the small sample size and complex, high-dimensional data characteristics. In summary, on this very challenging dataset, only PLURAL can produce reasonable prediction results.

## VI. CONCLUSION

In this paper, we presented PLURAL as an approach for transfer learning with unstructured, 3D point clouds. PLURAL contains a novel model architecture and a new contrastive learning framework, that leverages the physical dissimilarity across 3D point cloud datasets due to different sensor configurations and environments. Based on the observation that high-level shape features are more transferable than low-level geometry features of 3D shapes, we propose to integrate a domain-specific 3D encoder with a domain-agnostic 3D module. We then conducted contrastive instance alignment, which is augmented by hard sample mining. The experiments from simulation and case studies show the effectiveness of the proposed PLURAL framework. Future work in this direction could investigate other related learning tasks (e.g., weak supervision), and apply this methodology to more applications. Additionally, we noticed a higher prevalence of self-supervised approaches in point cloud transfer learning when contrasted with multi-domain approaches. This could stem from the challenges associated with uncertain cross-domain variance and uneven inter-class distribution, as highlighted by Huang et al. [47]. Nevertheless, conducting a thorough exploration into the distinct mechanisms through which various transfer learning methods handle domain-specific knowledge and exchange domain-invariant knowledge would constitute an interesting avenue for future research.

## REFERENCES

[1] S. R. Richter, V. Vineet, S. Roth, and V. Koltun, "Playing for data: Ground truth from computer games," in *Proc. Eur. Conf. Comput. Vis.*, Amsterdam, The Netherlands. Cham, Switzerland: Springer, Oct. 2016, pp. 102–118.

[2] M. Cordts et al., "The cityscapes dataset," in *Proc. CVPR Workshop Future Datasets Vis.*, vol. 2, 2015.

[3] A. Geiger, P. Lenz, C. Stiller, and R. Urtasun, "Vision meets robotics: The KITTI dataset," *Int. J. Robot. Res.*, vol. 32, no. 11, pp. 1231–1237, Sep. 2013.

[4] M. Imad, O. Doukhi, and D.-J. Lee, "Transfer learning based semantic segmentation for 3D object detection from point cloud," *Sensors*, vol. 21, no. 12, p. 3964, Jun. 2021.

[5] Z. Yan, L. Sun, T. Duckctr, and N. Bellotto, "Multisensor online transfer learning for 3D LiDAR-based human detection with a mobile robot," in *Proc. IEEE/RSJ Int. Conf. Intell. Robots Syst. (IROS)*, Oct. 2018, pp. 7635–7640.

[6] Z. Chai and C. Zhao, "Fault-prototypical adapted network for cross-domain industrial intelligent diagnosis," *IEEE Trans. Autom. Sci. Eng.*, vol. 19, no. 4, pp. 3649–3658, Oct. 2022.

[7] Z. Chai, C. Zhao, and B. Huang, "Multisource-refined transfer network for industrial fault diagnosis under domain and category inconsistencies," *IEEE Trans. Cybern.*, vol. 52, no. 9, pp. 9784–9796, Sep. 2022.

[8] R. Q. Charles, H. Su, M. Kaichun, and L. J. Guibas, "PointNet: Deep learning on point sets for 3D classification and segmentation," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, Jul. 2017, pp. 77–85.

[9] N. Audebert, B. Le Saux, and S. Lefèvre, "Semantic segmentation of Earth observation data using multimodal and multi-scale deep networks," in *Proc. 13th Asian Conf. Comput. Vis. Comput. Vis. (ACCV)*, Taipei, Taiwan. Cham, Switzerland: Springer, Nov. 2017, pp. 180–196.

[10] D. Maturana and S. Scherer, "VoxNet: A 3D convolutional neural network for real-time object recognition," in *Proc. IEEE/RSJ Int. Conf. Intell. Robots Syst. (IROS)*, Sep. 2015, pp. 922–928.

[11] C. R. Qi, L. Yi, H. Su, and L. J. Guibas, "PointNet++: Deep hierarchical feature learning on point sets in a metric space," in *Proc. Adv. Neural Inf. Process. Syst.*, vol. 30, 2017, pp. 1–10.

[12] W. Wu, Z. Qi, and L. Fuxin, "PointConv: Deep convolutional networks on 3D point clouds," in *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit. (CVPR)*, Jun. 2019, pp. 9613–9622.

[13] Q. Hu et al., "RandLA-Net: Efficient semantic segmentation of large-scale point clouds," in *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit. (CVPR)*, Jun. 2020, pp. 11105–11114.

[14] M.-H. Guo, J.-X. Cai, Z.-N. Liu, T.-J. Mu, R. R. Martin, and S.-M. Hu, "PCT: Point cloud transformer," *Comput. Vis. Media*, vol. 7, no. 2, pp. 187–199, Jun. 2021.

[15] J. Tobin, R. Fong, A. Ray, J. Schneider, W. Zaremba, and P. Abbeel, "Domain randomization for transferring deep neural networks from simulation to the real world," in *Proc. IEEE/RSJ Int. Conf. Intell. Robots Syst. (IROS)*, Sep. 2017, pp. 23–30.

[16] A. Pashevich, R. Strudel, I. Kalevatykh, I. Laptev, and C. Schmid, "Learning to augment synthetic images for Sim2Real policy transfer," in *Proc. IEEE/RSJ Int. Conf. Intell. Robots Syst. (IROS)*, Nov. 2019, pp. 2651–2657.

[17] N. Liu, Y. Cai, T. Lu, R. Wang, and S. Wang, "Real-sim-real transfer for real-world robot control policy learning with deep reinforcement learning," *Appl. Sci.*, vol. 10, no. 5, p. 1555, Feb. 2020.

[18] K. Arndt, M. Hazara, A. Ghadirzadeh, and V. Kyrki, "Meta reinforcement learning for sim-to-real domain adaptation," in *Proc. IEEE Int. Conf. Robot. Autom. (ICRA)*, May 2020, pp. 2725–2731.

[19] C. Wu, X. Bi, J. Pfrommer, A. Cebulla, S. Mangold, and J. Beyerer, "Sim2real transfer learning for point cloud segmentation: An industrial application case on autonomous disassembly," in *Proc. IEEE/CVF Winter Conf. Appl. Comput. Vis. (WACV)*, Jan. 2023, pp. 4520–4529.

[20] S. Horache, J.-E. Deschaud, and F. Goulette, "3D point cloud registration with multi-scale architecture and unsupervised transfer learning," in *Proc. Int. Conf. 3D Vis. (3DV)*, Dec. 2021, pp. 1351–1361.

[21] A. Xiao, J. Huang, D. Guan, F. Zhan, and S. Lu, "Transfer learning from synthetic to real LiDAR point cloud for semantic segmentation," in *Proc. AAAI Conf. Artif. Intell.*, 2022, vol. 36, no. 3, pp. 2795–2803.

[22] S. Xie, J. Gu, D. Guo, C. R. Qi, L. Guibas, and O. Litany, "Point-Contrast: Unsupervised pre-training for 3D point cloud understanding," in *Proc. Eur. Conf. Comput. Vis.*, Glasgow, U.K. Cham, Switzerland: Springer, Aug. 2020, pp. 574–591.

[23] J. Wei, G. Lin, K.-H. Yap, T.-Y. Hung, and L. Xie, "Multi-path region mining for weakly supervised 3D semantic segmentation on point clouds," in *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit. (CVPR)*, Jun. 2020, pp. 4383–4392.

[24] B. Du, X. Gao, W. Hu, and X. Li, "Self-contrastive learning with hard negative sampling for self-supervised point cloud learning," in *Proc. 29th ACM Int. Conf. Multimedia*, Oct. 2021, pp. 3133–3142.

[25] F. Yang, Y. Cao, Q. Xue, S. Jin, X. Li, and W. Zhang, "Contrastive embedding distribution refinement and entropy-aware attention for 3D point cloud classification," 2022, *arXiv:2201.11388*.

[26] Y. Zeng et al., "RT3D: Real-time 3-D vehicle detection in LiDAR point cloud for autonomous driving," *IEEE Robot. Autom. Lett.*, vol. 3, no. 4, pp. 3434–3440, Oct. 2018.

[27] A. Shrivastava, A. Gupta, and R. Girshick, "Training region-based object detectors with online hard example mining," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, Jun. 2016, pp. 761–769.

[28] R. Li, X. Li, P.-A. Heng, and C.-W. Fu, "PointAugment: An auto-augmentation framework for point cloud classification," in *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit. (CVPR)*, Jun. 2020, pp. 6377–6386.

[29] J. Snoek, H. Larochelle, and R. P. Adams, "Practical Bayesian optimization of machine learning algorithms," in *Proc. Adv. Neural Inf. Process. Syst.*, vol. 25, 2012, pp. 1–9.

[30] D. R. Jones, "A taxonomy of global optimization methods based on response surfaces," *J. Global Optim.*, vol. 21, no. 4, pp. 345–383, 2001.

[31] A. T. Bui and D. W. Apley, "Analyzing nonparametric part-to-part variation in surface point cloud data," *Technometrics*, vol. 64, no. 4, pp. 457–474, Oct. 2022.

[32] H. Yan, K. Paynabar, and M. Pacella, "Structured point cloud data analysis via regularized tensor regression for process modeling and optimization," *Technometrics*, vol. 61, no. 3, pp. 385–395, Jul. 2019.

[33] D. Girardeau-Montaut, M. Roux, R. Marc, and G. Thibault, "Change detection on points cloud data acquired with a ground laser scanner," *Int. Arch. Photogramm., Remote Sens. Spatial Inf. Sci.*, vol. 36, no. 3, p. W19, 2005.

[34] T. Uchida et al., "Noise-robust transparent visualization of large-scale point clouds acquired by laser scanning," *ISPRS J. Photogramm. Remote Sens.*, vol. 161, pp. 124–134, Mar. 2020.

[35] J. Wang and J. Shan, "Space filling curve based point clouds index," in *Proc. 8th Int. Conf. GeoComput.*, 2005, pp. 551–562.

[36] I. Rychkov, J. Brasington, and D. Vericat, "Computational and methodological aspects of terrestrial surface analysis based on point clouds," *Comput. Geosci.*, vol. 42, pp. 64–70, May 2012.

[37] H. Miao, A. Wang, B. Li, T.-S. Chang, and J. Shi, "Process modeling with multi-level categorical inputs via variable selection and level aggregation," *IISE Trans.*, vol. 55, no. 4, pp. 363–376, Apr. 2023.

[38] K. Dransfield, K. Addinall, and P. Bills, "Comparison and appraisal of techniques for the determination of material loss from tapered orthopaedic surfaces," *Wear*, vols. 478–479, Aug. 2021, Art. no. 203903.

[39] G. Moroni and S. Petro, "Geometric tolerance evaluation: A discussion on minimum zone fitting algorithms," *Precis. Eng.*, vol. 32, no. 3, pp. 232–237, Jul. 2008.

[40] M. Biehler, H. Yan, and J. Shi, "ANTLER: Bayesian nonlinear tensor learning and modeler for unstructured, varying-size point cloud data," *IEEE Trans. Autom. Sci. Eng.*, pp. 1–14, Jan. 2023, doi: 10.1109/TASE.2022.3230563.

[41] B. Sun, J. Feng, and K. Saenko, "Correlation alignment for unsupervised domain adaptation," in *Domain Adaptation in Computer Vision Applications*. Cham, Switzerland: Springer, 2017, pp. 153–171.

[42] V. Vincent, M. Wannes, and D. Jesse, "Transfer learning for anomaly detection through localized and unsupervised instance selection," in *Proc. AAAI Conf. Artif. Intell.*, 2020, vol. 34, no. 4, pp. 6054–6061.

[43] S. J. Pan, I. W. Tsang, J. T. Kwok, and Q. Yang, "Domain adaptation via transfer component analysis," *IEEE Trans. Neural Netw.*, vol. 22, no. 2, pp. 199–210, Nov. 2010.

[44] G. Wu, T. Cao, B. Liu, X. Chen, and Y. Ren, "Towards universal LiDAR-based 3D object detection by multi-domain knowledge transfer," in *Proc. IEEE/CVF Int. Conf. Comput. Vis.*, Jun. 2023, pp. 8669–8678.

[45] C. Qin, H. You, L. Wang, C.-C. J. Kuo, and Y. Fu, "PointDAN: A multi-scale 3D domain adaption network for point cloud representation," in *Proc. Adv. Neural Inf. Process. Syst.*, vol. 32, 2019, pp. 1–12.

[46] J. Haurum, M. Allahham, M. Lynge, K. Henriksen, I. Nikolov, and T. Moeslund, "Sewer defect classification using synthetic point clouds," in *Proc. 16th Int. Joint Conf. Comput. Vis., Imag. Comput. Graph. Theory Appl.*, 2021, pp. 891–900.

[47] S. Huang, B. Zhang, B. Shi, H. Li, Y. Li, and P. Gao, "SUG: Single-dataset unified generalization for 3D point cloud classification," in *Proc. 31st ACM Int. Conf. Multimedia*, Oct. 2023, pp. 8644–8652.

**Michael Biehler** received the B.S. and M.S. degrees in industrial engineering with a major in production engineering from the Karlsruhe Institute of Technology (KIT) in 2017 and 2020, respectively. He is currently pursuing the Ph.D. degree with the H. Milton Stewart School of Industrial and Systems Engineering, Georgia Institute of Technology. His research interests include the interface between 3D machine learning and cyber-physical security, with the goal of developing methods for efficient and safe decision-making in complex (manufacturing) systems.

**Yiqi Sun** is currently pursuing the degree with the College of Computing, Georgia Institute of Technology. His long-term research goal is to build machines that can continually learn concepts from their experiences and apply them for reasoning and planning in the physical world.

**Shriyanshu Kode** is currently pursuing the degree in computer science with the Georgia Institute of Technology, with concentrations in intelligence and devices. He is interested in software development and machine learning from a wide range of applications, such as robotics and finance.

**Jing Li** received the Ph.D. degree in industrial and operations engineering from the University of Michigan, Ann Arbor, MI, USA. She is currently the Virginia C. and Joseph C. Mello Chair and a Professor with the H. Milton Stewart School of Industrial and Systems Engineering, Georgia Institute of Technology, Atlanta, GA, USA. Her research interests include statistical modeling and machine learning for healthcare applications. She is a member of IISE and INFORMS. She was a recipient of the NSF CAREER Award.

**Jianjun Shi** received the B.S. and M.S. degrees in automation from the Beijing Institute of Technology in 1984 and 1987, respectively, and the Ph.D. degree in mechanical engineering from the University of Michigan in 1992. Currently, he is the Carolyn J. Stewart Chair and a Professor with the H. Milton Stewart School of Industrial and Systems Engineering, Georgia Institute of Technology. His research interests include the fusion of advanced statistical and domain knowledge to develop methodologies for modeling, monitoring, diagnosis, and control of complex manufacturing systems. He is a fellow of four professional societies, including ASME, IISE, INFORMS, and SME, an Elected Member of the International Statistics Institute (ISI), a Life Member of ASA, an Academician of the International Academy for Quality (IAQ), and a member of the National Academy of Engineers (NAE).