NGTTA: NON-PARAMETRIC GEOMETRY-DRIVEN TEST TIME ADAPTION FOR 3D POINT CLOUD SEGMENTA TION

Anonymous authors

Paper under double-blind review

ABSTRACT

Previous Test Time Adaption (TTA) methods usually suffer from training collapse when they are transferred to complex 3D scenes for point cloud segmentation due to the significant domain gap between the source and target data. To solve this issue, we propose NGTTA, a stable test time adaption method guided by non-parametric geometric features. In NGTTA, we leverage the distribution of non-parametric geometric features on target data as an "intermediate domain" to reduce the domain gap and guide the stable learning of the source model on target data. Specifically, we use the source domain model and a non-parametric geometric model to extract the embedding features and geometric features of the point cloud, respectively. Then, a category-balance sampler is designed to filter easy samples and hard samples in the input data to address the class imbalance issue in semantic segmentation. Inspired by previous work, we use easy samples for entropy minimization loss and pseudo-label prediction to fine-tune the source domain model. The difference is that we refine the pseudo labels not only by considering the soft voting among their nearest neighbors in the model embedding feature space but also in the geometric space, which can prevent the accumulation of errors caused by model feature shifts. Furthermore, we believe that hard samples can effectively represent the distribution differences between the source domain and the target domain. Therefore, we propose to distill the geometric features of hard samples into the source domain model in the early stages of training to quickly converge to an "intermediate domain" that is similar to the target domain. By taking advantage of the ability of the non-parametric geometric feature to represent the underlying manifolds of the target data, our method efficiently reduces the difficulty of the domain adaption. We conduct the main experiments on the more challenge sim-to-real benchmark about synthetic dataset 3DFRONT and the realworld datasets ScanNet and S3DIS for 3D segmentation task. Results show that our method can efficiently improve the mIOU by over 3% on $3DFRONT \rightarrow ScanNet$ and 7% on 3DFRONT \rightarrow S3DIS.

041

006

008 009 010

011

013

014

015

016

017

018

019

021

024

025

026

027

028

029

031

032

034

1 INTRODUCTION

042 With the development of deep learning, more and more neural networks are deployed in real-world 043 applications. However, current deep networks may only perform optimally when the training and 044 testing data share the same distribution (He et al., 2016; Krizhevsky et al., 2017). Therefore, deep networks often struggle to generalize on the unseen data which is known as the domain shift (Geirhos et al., 2018; Hendrycks & Dietterich, 2019; Recht et al., 2019). Unsupervised Domain 046 Adaptation (UDA) techniques have emerged as a popular solution for addressing domain shift in deep 047 learning (Saito et al., 2017; Peng et al., 2021; Wu et al., 2019; Long et al., 2017). These methods aim 048 to transfer knowledge from a labeled source domain to an unlabeled target domain during training. While UDA methods have demonstrated effectiveness in improving the performance of deep networks in the presence of domain shift, a key limitation is the requirement to have knowledge of the test data 051 during training. This constraint can greatly diminish the practical utility of UDA techniques. 052

Recently, an interesting and practical paradigm known as Test Time Adaption (TTA) is attracting more and more attention (Wang et al., 2021; Chen et al., 2022; Niu et al., 2023; Zhang et al., 2022a).



Figure 1: Ilustration of the overall pipeline. i) In Stage 1, the model was trained on the source domain
with the given label. From the source distribution, we can see that the class distribution is better
separated from each other due to the correct supervision signal. ii) In Stage 2, due to the significant
differences in feature distribution, transferring directly from S2T-Distribution to Target-Distribution
based on previous TTA methods is difficult. (For example, class distribution inside the yellow box in
Target Distribution overlaps almost completely in S2T-Distribution). iii) We use the non-parametric
geometric feature, which can be considered as the 'intermediate domain' between the source domain
and the target domain to guide the adaptation process.

071 It does not need to access the training method and training data used by the model and can adapt 072 any trained source model to the test data in testing time. This flexibility and independence make 073 TTA a valuable approach for addressing domain shift in real-world scenarios. However, previous 074 TTA methods often perform training collapse on the 3D segmentation task, especially on the more 075 challenging sim-to-real benchmark proposed in (Ding et al., 2022). By digging into the failure cases, 076 we found the two main challenges: i) Previous TTA methods often use entropy minimization loss on classification tasks, but when this loss is applied to segmentation tasks, it may cause the model to be 077 overconfident in the majority class due to the more serious class imbalance problem. ii) 3D scenes exhibit higher complexity compared to images, and the distribution variations among different scenes 079 are significantly greater, so using a simple TTA method to complete this difficult process may lead to the training collapse, which is shown in Figure 1. 081

The first challenge arises mainly because previous TTA methods set a fixed low threshold to select low-entropy easy samples for entropy minimization loss. However, due to the common issue of class imbalance in semantic segmentation, low-entropy samples are predominantly found in the major categories, which exacerbates the class imbalance problem further. To address this issue, we propose a category-balance sampler. Unlike previous methods that set a fixed threshold, we calculate a corresponding entropy threshold for each category based on its sample count. This approach reduces the differences in the number of easy samples across categories, effectively mitigating the class imbalance problem. Then, we propose to use easy samples for entropy minimization loss and pseudo-label prediction loss.

To solve the second challenge, we hope to find an "intermediate domain" between the source domain 091 and target domain to guide the adapting process of the model thereby promoting training stability. 092 Inspired by (Ran et al., 2022; Sun et al., 2024), explicit geometric representations can capture the underlying manifolds of the data, thereby offering insights into the rough distribution of the target 094 domain to improve the model's generalization to unseen data. Therefore we use the non-parametric 095 geometric feature as the "intermediate domain" between the source domain to the target domain 096 which is shown in Figure 1. Specifically, we additionally utilize a non-parametric geometric model to obtain the geometric features of the point cloud. First, for the easy samples, we recognize that 098 the source domain model has a feature-shifting problem on the unseen data, which may lead to incorrect pseudo-labels and subsequently cause the accumulation of errors. Therefore, we refine the 099 pseudo labels not only by considering the soft voting among their nearest neighbors in the model's 100 embedding feature space but also in the geometric space. Then, we believe that hard samples can 101 effectively represent the differences between the source domain and the target domain. Therefore, we 102 distill the geometric features of the hard samples into the source domain model. This process helps 103 the source domain model quickly learn the underlying manifold distribution of the target domain, 104 converging to an "intermediate domain" that closely resembles the target domain. This approach 105 enhances the stability and performance of domain adaptation. 106

107 Testing time adaptation for indoor 3D scene segmentation tasks remains an unexplored area. To the best of our knowledge, we are the first to attempt such an approach. Therefore, we propose to follow

the UDA (Unsupervised Domain Adaptation) method Ding et al. (2022) and conduct experiments using the challenging sim2real benchmark. Our method can efficiently improve the source model by over 3% mIOU on 3DFRONT \rightarrow ScanNet and 7% on 3DFRONT \rightarrow S3DIS. Our contributions can be summarized as follows:

- We proposed a category-balance sampler to filter easy samples and hard samples, ensuring that the difference in the number of positive samples for each category is minimized, thereby addressing the class imbalance problem.
- We propose to leverage the non-parametric geometric feature as the "intermediate domain" to stabilize the adaption process to rapidly converge the model to a distribution that approximates the target domain.
 - We conducted experiments on the challenging sim2real benchmark, and the competitive experimental results validate the effectiveness of our method.
- 120 121

112

113

114

115

116

117

118

119

2 RELATED WORK

122 123

124 **Unsupervised Domain Adaptation:** UDA (Saito et al., 2017; Peng et al., 2021; Wu et al., 2019; 125 Long et al., 2017; Yang et al., 2020; Zou et al., 2018; Cui et al., 2020) aims at transferring knowledge 126 of the labeled source domain to the unlabeled target domain in the training time. Segmentation tasks 127 are more difficult to perform domain generalization than simple classification tasks, especially on the 3D data. Zou et al. (2018) improved the performance by solving the class imbalance problem in the 128 segmentation task. Squeezesegv2 (Wu et al., 2019) consider the density and geometric during the 129 domain adaption. Jaritz et al. (2020) propose to leverage the information of images and point cloud 130 to complete multi-modality UDA. 131

132 **Test Time Adaption:** Almost all previous TTA Methods have been applied to classification tasks. Wang et al. (2021) firstly proposes fully test time adaption which does not need to access the training 133 method and training data and uses entroy minimization loss to optimize the model. Then, many 134 subsequent works (Zhao et al., 2023; Niu et al., 2023; Zhang et al., 2022a; Wang et al., 2022; Niu 135 et al., 2022) modify the entropy loss to further improve the performance. Chen et al. (2022) leverages 136 self-supervised contrastive learning and a soft voting strategy for refining the pseudo-label to facilitate 137 target feature learning. Iwasawa & Matsuo (2021) tries to update the category prototypes on the 138 target domain to provide a more accurate decision boundary. Howevere, there have been few test 139 time adaption methods focusing on the segmentation task, especially on the more difficult 3D data. 140 Song et al. (2023) firstly explore TTA method for segmentation in the dynamic world. Shin et al. 141 (2022) propose a multi-modal test time adaption framework for 3D segmentation. However, the need 142 for multi-modal limits its general applications. Therefore, it is necessary to develop a general and 143 effective TTA method for 3D segmentation.

144 **3D Point Cloud Segmentation:** Due to the disorder and irregularity of the point cloud data, Qi 145 et al. (2017) firstly propose PointNet++ to use ball query or kNN to construct local neighborhood 146 and aggregate it by symmetry pooling. Many subsequent works (Qian et al., 2022; Lin et al., 147 2023; Thomas et al., 2019; Zhao et al., 2021) design more complex modules to extract the local 148 feature based on the PointNet++. PointTransformer (Zhao et al., 2021) leverage the local attention 149 mechanism to extract local feature. KPConv (Thomas et al., 2019) defined the anchor points and used them to compute the aggregate weight. In addition, Ran et al. (2022); Sun et al. (2024) propose to 150 leverage the explicit geometric to introduce strong prior which can reduce the learning difficulty and 151 improve the performance. Furthermore, PointNN (Zhang et al., 2023) uses trigonometric functions to 152 capture the non-parametric geometric feature for point cloud recognition which has proved the strong 153 generalization ability to unseen 3D data. Inspired by the above methods, we decided to leverage the 154 non-parametric geometric feature to prompt the learning of test time adaption.

155 156

3 PROPOSED METHOD

157 158

We address the closed-set fully test time adaption in 3D segmentation task that we can only access the source model during the adaption process. As shown in Figure 1, the source model is trained on the labeled source domain $\{x_i^s, y_i^s\}_{i=1}^{N_s}$, where x_i^s is the input *i*-th scene in source dataset, y_i^s is the corresponding label, and N_s is the total number of scenes in source dataset. The goal of



Figure 2: Ilustration of the proposed framework. We use a Category-Balance Sampler to balance the selection of simple and hard samples for each category. Simple samples are used for entropy minimization loss and pseudo-label prediction loss, with pseudo-labels optimized through soft voting using features from the non-parameterized geometric model and the source domain model. hard samples are used to distill their geometric features into the source domain model.

test time adaption is to adapt the source model on the target domain $\{x_i^t\}_{i=1}^{N_t}$ without accessing its labels $\{y_i^t\}_{i=1}^{N_t}$ during the adaption process, where N_t is the total number of scenes in the target dataset. It's worth noting that our setting is closed-set which means that both source domain and target domain share the same semantic classes. Therefore, there will be $y_i^s = 0, 1, ..., N_c - 1$ and $y_i^s = 0, 1, ..., N_c - 1$, where N_c is the total number of classes. Since test time adaption only considers training on the target domain, we use x_i and y_i to replace the x_i^t and y_i^t for convenience and the following symbols are all defined on the target domain.

The framework of NGTTA is shown in Figure 2. Firstly, We designed a category-balance sampler to 189 balance the selection of easy samples and difficult samples and use easy samples to perform entropy 190 minimization loss and pseudo-label prediction loss, which will be introduced in 3.1. Then, in 3.2 we 191 additionally propose to use a non-parametric geometric model to extract the geometric information of 192 point clouds. For easy samples, we refine the pseudo-labels by modified soft voting, which averages 193 the neighborhood predictions from both the model feature space and the geometric feature space. For 194 difficult samples, we distill their geometric features into the source domain model, helping the source 195 domain model quickly converge to an "intermediate domain" that is similar to the target domain. 196

3.1 CATEGORY-BALANCE SAMPLER

Our method is generally applicable to point-based models. Therefore, for any model \mathcal{M}^s that is trained on the source domain, the output for the input scene $x_i \in \mathbb{R}^{M \times 3}$ in the target domain will be the point-wide embedding feature $f_i \in \mathbb{R}^{M \times C}$, where M represents the number of points in *i*-th scene and C represents the number of feature channel.

$$f_i = \mathcal{M}^s(x_i) \tag{1}$$

Then, f_i is sent to the classification head to produce the class probability $p_i \in R^{M \times N_c}$.

$$p_i = Head(f_i) \tag{2}$$

We calculate the entropy E_i of the samples based on the class probability p_i , which can represent the confidence level of the model's predictions.

211

197

198

203

204 205

206 207 208

175

176

177

178

179

181

213

 $E_{i} = -\sum_{c=0}^{N_{c}-1} p_{i}[c] \log p_{i}[c]$ (3)

215 Next, we need to determine whether each sample is an easy sample or a hard sample based on its entropy and the corresponding threshold. Previous methods shared the same threshold for samples

across all categories to select easy low-entropy samples. However, indoor scene segmentation
 generally suffers from class imbalance issues, where the majority of confidently predicted low entropy samples are found in the major categories. As a result, samples from tail categories are often
 difficult to select, which further exacerbates the class imbalance problem.

Therefore, we propose a category-balance sampler. Specifically, we first calculate the number of samples Z_c for the *c*-th category and define the category with the highest sample count as the major category, with the count being Z_m . Then, we define the threshold for *c*-th category as follows:

$$\sigma_c = \sigma + (1 - \frac{Z_c}{Z_m})\gamma\tag{4}$$

where σ is the initial threshold, and γ adaptively adjusts the threshold based on the number of samples in each category. We can see that as the sample count decreases, the threshold will gradually increase. This means that tail categories with fewer samples will have a larger entropy threshold, allowing for the selection of more samples to reduce the disparity between the number of samples across different categories.

Then, we define the set of easy samples \mathcal{G}_i^e as follows:

224

225

233

239 240

244 245

254 255 256

257 258

259

$$\mathcal{G}_i^e = \{j \mid E_{ij} < \sigma_{y'_{ij}}\} \tag{5}$$

where j means the j-th sample in the i-th scene. y'_{ij} represents the class prediction of the j-th sample, where $y'_{ij} = Argmax(p_{ij})$.

In contrast, the set of hard samples \mathcal{G}_i^h is defined as follows:

$$\mathcal{G}_i^h = \{ j \mid E_{ij} \ge \sigma_{y'_{ii}} \} \tag{6}$$

Inspired by previous TTA methods Wang et al. (2021), we apply entropy minimization loss to low-entropy samples to avoid increasing the confidence level of incorrect predictions, which can be written as follows:

$$L_{ent} = \underset{j \in \mathcal{G}_i^e}{Min(E_{ij})} \tag{7}$$

246 However, entropy minimization loss is category-agnostic and can only enhance the model's confidence 247 level. To improve the model's performance on semantic segmentation metrics, we introduce pseudolabel prediction loss. Simply put, we copy the source domain model as a momentum model to predict 248 pseudo-labels. Unlike the source domain model, we fix the parameters of the momentum model 249 during training, and every n epochs, we copy the parameters from the source domain model to the 250 momentum model to ensure the stability of the pseudo-labels. We then use the class predictions from 251 the source domain model together with the pseudo-labels to compute the classification loss. The 252 process can be written as follows: 253

$$L_{cls} = CrossEntropy(p_{ij}, y_{ij}^m)$$

$$(8)$$

where y_{ij}^m means the pseudo-label from the momentum model.

3.2 NON-PARAMETRIC GEOMETRY-DRIVEN ADAPTION

260 As discussed above, due to the significant domain gap between source and target in 3D segmentation 261 task, the adaption process from source to target is difficult. Some domain adaption methods (Li et al., 2021; Wang et al., 2023) propose to define an "intermediate domain" that guides the source model to 262 progressively adapt to the target domain. However, those methods need to access the source data, 263 which is not suitable for our setting. Inspired by (Ran et al., 2022; Sun et al., 2024), non-parametric 264 geometric feature can capture the underlying manifolds of point cloud data which can represent the 265 rough approximation of the target distribution, so we leverage non-parametric geometric feature as 266 the "intermediate domain" to boost the test time adaption for 3D segmentation. 267

We leverage PointNN (Zhang et al., 2023) as the non-parametric model which has a strong general ization ability from seen to unseen 3D data. For PointNN, it uses farthest point sampling and kNN to downsample and construct the neighborhood like PointNet++ (Qi et al., 2017). However, PointNet++

uses learnable MLP to extract neighborhood point features, while PointNN uses non-parametric
 trigonometric functions. Since it does not require training, PointNN can be directly applied to
 unlabeled test time adaption process.

273 274 Specifically, for the input scene x_i , non-parametric model \mathcal{M}^g output the geometric feature f_i^g , which can be written as follows:

$$f_i^g = \mathcal{M}^g(x_i) \tag{9}$$

Easy Feature Soft Voting. Although easy samples have a higher confidence level, there are still
 incorrect predictions that lead to the provision of erroneous pseudo-labels. To address this issue,
 we propose constructing K-nearest neighbors using the features of easy samples and correcting the
 pseudo-labels through soft voting based on the neighbor class predictions.

However, due to the feature shift phenomenon of the source domain model in the target domain, the K-nearest neighbors may not be accurate, which affects the correction of the pseudo-labels. To address this issue, we simultaneously consider K-nearest neighbors constructed using geometric features during the soft voting process. Specifically, for the *j*-th easy sample feature of source model and the non-parametric model f_{ij} and f_{ij}^g , where $j \in \mathcal{G}_i^e$. We construct the neighborhood Q_{ij} and Q_{ij}^g by kNN. Then we use the y_i^m from momentum model to generate the new pseudo-label, which can be written as follows:

$$Y_{ij}^{ps} = \beta_t \frac{1}{K} \sum_{k \in Q_{ij}} y_{ik}^m + (1 - \beta_t) \frac{1}{K} \sum_{k \in Q_{ij}^g} y_{ik}^m$$
(10)

where Y_{ij}^{ps} is the refined pseudo-label of the *j*-th sample in *i*-th scene. And the β_t is the weight factor that gradually changes as training progresses. In simple terms, we believe that as the source domain model gradually converges to the target domain, the feature shift phenomenon decreases, and thus the importance of geometric features will diminish. Therefore, we assign higher weight to geometric features in the early stages of training, while in the later stages, we assign higher weight to the model's embedding features.

$$\beta_t = \frac{t}{T}\beta \tag{11}$$

where t is the current training step and T is the total number of training steps.

Finally, we modify Eq. 8 as follows to achieve better classification loss. $I_{L} = C_{\text{magas}} E_{\text{retromy}}(n - V_{ps}^{ps})$

$$L_{cls} = CrossEntropy(p_{ij}, Y_{ij}^{ps})$$
(12)

Hard Feature Distillation. Previous methods typically consider hard samples as noisy samples, thus
 only handling easy samples while discarding hard ones. However, we believe that many hard samples
 arise from the significant differences between the target domain and the source domain in the context
 of scene segmentation sim2real benchmarks, and therefore they can represent the information about
 the distributional differences. To this end, we decide to utilize this difference information to facilitate
 the adaptation of the source domain model to the target domain. Rather than using the erroneous
 predictions of hard samples for the aforementioned two losses, we believe that distilling their features
 to transfer distributional information is a more effective approach.

Therefore, we propose a feature distillation loss aimed at distilling geometric features into the source domain model, enabling it to quickly converge to an "intermediate domain" that is closer to the target domain, thereby stabilizing the adaptation process, which can be written as:

$$L_{dis} = \sum_{j \in \mathcal{G}_i^h} MSE(MLP(f_{ij}), f_{ij}^g)$$
(13)

where MLP is a multil-ayer perceptron used to align the dimensions of source model feature with the geometric feature.

4 EXPERIMENTS

319 320

322

314 315

316

317 318

276

288 289

296

297

301

321 4.1 DATASETS

We conducted experiments on three datasets consists of a synthetic dataset 3DFRONT and the real-world datasets ScanNet and S3DIS for 3D segmentation task.

Table 1: Test Time Adaption Results of Sim-to-Real (3DFRONT \rightarrow ScanNet and 3DFRONT \rightarrow S3DIS) Benchmark . We report the mIOU (%). mACC (%) and OA(%) of different UDA and TTA methods. **Bold** represents the best performance in UDA and TTA methods

Tuna	Mathad	3DFF	RONT→So	anNet	3DFI	RONT→S	3DIS
Type	Method	mIOU	mACC	OA	mIOU	mACC	OA
	Source Only	34.80	49.24	71.66	29.38	39.72	68.47
	SqueezeSegV2 (Wu et al., 2019)	34.98	49.72	71.89	29.80	40.12	69.81
UDA	AdaptSegNet (Tsai et al., 2018)	40.23	52.16	75.10	37.12	49.82	76.23
	APO-DA (Yang et al., 2020)	37.82	50.92	73.27	35.21	47.69	74.91
	TENT (Wang et al., 2021)	15.63	28.62	55.41	31.98	42.38	71.23
TTA	DOT (Zhao et al., 2023)	18.30	29.71	56.58	32.61	43.41	72.30
IIA	AdaContrast (Chen et al., 2022)	30.57	49.61	73.05	33.28	45.20	73.01
	MEMO (Zhang et al., 2022a)	14.21	27.93	54.96	27.10	37.98	67.10
	T3A (Iwasawa & Matsuo, 2021)	17.20	29.17	56.02	32.11	43.11	71.80
TTA	Ours	38.42	51.30	74.03	36.36	49.07	75.74

ScanNet is proposed in (Dai et al., 2017), which is a popular real-world 3D scene dataset with 1,201
 scans for training, 3,12 scans for validation and 100 scans for testing. It has rich dense segmentation
 annotations for 20 categories.

S3DIS is proposed in (Armeni et al., 2016), which is a real-world 3D scene dataset with 271 scenes and rich dense segmentation annotations for 13 categories. Following the previous work (Qi et al., 2017), we used Area5 as validation set and others as training set.

348
 349
 350
 350
 351
 352
 352
 353
 354
 355
 356
 357
 358
 359
 359
 350
 351
 351
 352
 353
 353
 354
 355
 355
 356
 357
 358
 359
 359
 350
 350
 351
 351
 352
 351
 352
 353
 354
 355
 355
 355
 356
 357
 357
 358
 358
 359
 350
 351
 351
 351
 352
 351
 352
 351
 352
 351
 352
 351
 352
 351
 352
 352
 351
 352
 351
 352
 352
 353
 354
 355
 355
 355
 355
 355
 356
 357
 357
 358
 358
 350
 351
 351
 352
 352
 353
 354
 355
 355
 355
 355
 356
 357
 357
 358
 358
 359
 350
 350
 351
 351
 351
 352
 352
 353
 354
 355
 355
 356
 356
 357
 357
 358
 358
 350
 350
 351
 351
 351
 351

Closed-Set Setting. Test Time Adaption is the closed-set setting that the source and the target domain share the same categories. Therefore, we follow the setting in (Ding et al., 2022) to select 11 categories for 3DFRONT→ ScanNet and 3DFRONT→ S3DIS. To better demonstrate the generality of our approach, we also select 8 categories for the domain adaption between ScanNet and S3DIS. The selected categories are all the categories that are shared between the two datasets..

4.2 IMPLEMENTATION

BackBone. To prove the effectiveness of our method, we used the state-of-the-art model Point Meta (Lin et al., 2023) as the source model by default in the following experiments. We also reported the performance of other models to demonstrate the applicability of our method.

364 **Optimizable Parameters.** How to determine the optimal parameters is important in test time adaption. Previous research denotes that the knowledge of data domain is saved in "BatchNorm". 366 Therefore, TENT propose to only update the BatchNorm Parameters. However, due to the complexity 367 of 3D data, scenarios may not be independently and identically distributed among themselves, which 368 greatly affects the performance of optimizing BatchNorm. AdaContrast (Chen et al., 2022) use 369 contrast learning to optimize the entire model. However, in 3D segmentation task, the model is also 370 more complex than the classification model, which is often the "Encoder-Decoder" architecture that 371 the optimization is very difficult. Therefore, in our implementation, we only optimized the embedding layer at the beginning of the model and the classification head at the end. 372

373

359

360

374 4.3 RESULTS375

Sim-to-Real Benchmark We tested different UDA and TTA methods on this benchmark which is
 shown in Table 1. We see that due to the significant domain gap between the synthetic and real-world
 datasets, the performance of the source model is only about 30% mIOU. We first report the current

7

325

326

327 328

Tuna	Mathad	S3I	DIS→Scan	Net	ScanNet→S3DIS			
Type	Method	mIOU	mACC	OA	mIOU	mACC	OA	
	Source Only	48.20	65.01	78.18	50.99	61.65	77.66	
	SqueezeSegV2 (Wu et al., 2019)	46.31	63.17	76.27	51.20	62.12	77.93	
UDA	AdaptSegNet (Tsai et al., 2018)	50.34	66.12	79.20	50.12	60.89	76.23	
	APO-DA (Yang et al., 2020)	53.35	68.87	81.92	52.65	63.01	78.62	
	TENT (Wang et al., 2021)	51.49	67.01	80.22	50.21	61.01	77.23	
TTA	DOT (Zhao et al., 2023)	51.92	67.51	80.67	51.00	61.72	77.92	
IIA	AdaContrast (Chen et al., 2022)	52.12	67.98	81.21	52.35	62.71	78.01	
	MEMO (Zhang et al., 2022a)	47.23	64.65	77.12	49.97	60.92	76.92	
	T3A (Iwasawa & Matsuo, 2021)	51.01	66.99	79.83	50.61	61.21	77.68	
TTA	Ours	52.71	68.23	81.62	53.78	64.25	79.34	

Table 2: Test Time Adaption Results of Cross-site Benchmark . We report the mIOU (%) of different UDA and TTA methods.

UDA methods, because the target domain data is accessed during training, the UDA method can effectively improve the performance, especially AdaptSegNet (Tsai et al., 2018) improving the due to the complexity of 3D scenes, as well as the sim-to-real difficulty, the TTA method does not perform optimally. For example, the miou of the most classical TENT method decreases on the 3DFRONT-ScanNet respectively about 19.17%. Afterward, through the experimental results, we find that DOT (Zhao et al., 2023) and AdaContrast (Chen et al., 2022) generally have better performance, because the former considers the class imbalance problem, and the latter introduces contrastive learning and pseudo-label adjustment strategies. Furthermore, we find that the TTA method generally performs better on the simpler 3DFRONT \rightarrow S3DIS. That is because the scene of S3DIS is more complete and easy than ScanNet and the domain gap with 3DFRONT is smaller. According to the results of previous methods, we believe that reducing the domain gap and solving the class imbalance is an effective way to improve the performance of TTA. Therefore, by introducing non-parametric geometric feature as "intermediate domain" and the class-balance entropy minimization loss, our method can improve the mIOU about 3.62% and 6.98% of 3DFRONT-ScanNet and UDA method.

411
412
413
414
414
414
414
414
414
414
414
414
415
415
416
416
417
417
418
419
419
419
410
410
410
411
411
411
412
412
413
414
414
415
415
416
416
417
417
418
419
419
419
410
410
411
411
411
411
412
412
412
413
414
414
415
415
415
416
417
417
418
419
419
410
411
411
411
411
412
412
412
413
414
415
415
415
416
417
417
417
418
419
419
419
410
410
411
411
411
412
412
412
413
414
415
415
415
416
417
417
418
419
419
419
410
410
411
411
411
411
412
412
412
412
413
414
414
415
415
416
417
417
418
418
419
419
419
419
410
410
410
411
411
411
412
412
412
413
414
414
415
415
416
417
416
417
418
418
419
419
419
410
410
410
411
411

Table 3: Test Time Adaption Results of Different Backbones. We use our method on different models of the 3DFRONT→S3DIS. Table 4: Comparison with pre-trained models.**NP** means non-parametric model and **PT** means pretrained model.

Models	mIOU	GFLOPs	TP	type	Pretrain Dataset	Models	mI
PointMetaBase-L (Lin et al., 2023) +ours	29.38 36.36	2.0 2.9	192 178	NP	-	PointNN	36
PointNet++ (Qi et al., 2017) +ours	22.74 28.10	7.2 7.9	181 171		ScanNet	PointM2AE CSC	34 35
PointTransformer (Zhao et al., 2021)	25.32	2.80	170	РТ		MSC	35
PointNeXt-L (Thomas et al., 2019) +ours	27.15	15.2 16.1	126 118		Structure3D	CSC MSC	36 36 36

430 More BackBones and Efficiency. To prove the generality of our method, we conducted the experiments of different backbones on the 3DFRONT→S3DIS which is shown in Table 3. We tested four backbones which consist of PointMeta, PointNet++, PointTransformer and PointNeXt. The

432 results show that our method is applicable to most point-based models and can effectively enhance 433 their performance, demonstrating the versatility of our approach. In addition, we also report the 434 efficiency of NGTTA, including the computational cost in GFLOPs and the inference speed measured 435 in Throughput (ins./sec). From the results, it can be observed that the additional computational 436 overhead we introduced is acceptable compared to the original model's expenses. This is mainly attributed to two design features: 1) We modified the original PointNN to reduce the calculation 437 cost by reducing the number of neighborhood points, feature dimensions and layers, which will be 438 introduced in supplementary material. 2) In the computationally intensive soft voting component, we 439 only selected a small proportion of clean samples to perform the operation, significantly reducing the 440 computational overhead. 441

442 443

4.4 ABLATION STUDY

444 **Comparison with pre-trained models** In this section, we compare the performance of using 445 parameterized pre-trained models (PointM2AE Zhang et al. (2022b), CSC Hou et al. (2021), MSC Wu 446 et al. (2023)) and non-parametric models to drive domain adaptation, which can be seen in Table 4. 447 The results indicate that the performance of the non-parametric model is superior. We believe this 448 is primarily due to two reasons: 1) The high-dimensional feature representation of parameterized 449 models is more abstract, which makes it difficult to facilitate the rapid convergence of the source 450 domain model. 2) Pre-trained models still learn information specific to the dataset, which can affect 451 the domain adaptation process. As the scale of the pre-training data increases, this influence gradually diminishes. For example, models pre-trained on Structure3D perform better than those pre-trained on 452 ScanNet. Therefore, compared to pre-trained models, non-parametric models can more explicitly 453 represent the features of the current data. Additionally, they do not contain any specific dataset 454 information and require no training steps. Thus, at the current stage, non-parametric models still 455 outperform parameterized pre-trained models. 456

457 Different Parts in Our Method. We tested the different parts in our method on the $3DFRONT \rightarrow S3DIS$ benchmark, which is shown in Table 5. There are three main parts in our 458 method, which consists of Distill (distillation from the non-parametric model), Soft Voting and 459 Category-Balance Sampler. From the results we see that all three modules contribute to the perfor-460 mance improvement, of which Distill has the most significant improvement due to the significant 461 reduction of the domain gap by utilizing "intermediate domains". Soft Voting and Category-Balance 462 Sampler can effectively improve the performance by more than 1% mIOU because they provide more 463 accurate pseudo-labels and solve the class imbalance problem. 464

Table 5: Ablation study result of different parts in our method. We use PointMetaBase-468 469

Table 6: Ablation Study of Soft Voting on the $3DFRONT \rightarrow S3DIS$ benchmark.

L as t	he source	model and test it on t	the	Distance	Neighbor Numbers	mIOU
Distill	Soft Voting	Category-Balance Sampler mIOU		Feat	$\begin{vmatrix} 2\\10 \end{vmatrix}$	35.12 35.83
		29.38 32.28	8 8	Tout	20	35.01
\checkmark	\checkmark	34.92 31.76	2		2	35.43
✓	\checkmark	√ 30.62 ✓ 36.3 0	2 6	Feat+Geo	20	36.36 36.21

476 **Soft Voting.** We conducted the ablation study of soft voting in the Table 6. First, we tested soft 477 voting for building neighborhoods with only the features of the model (Feat) which is proposed in AdaContrast (Chen et al., 2022). We see that when the number of neighborhoods is 10, the 478 performance is the best, which also proves the importance of the soft voting strategy. However, 479 when the number of neighborhoods continues to increase, the offset model features may introduce 480 neighborhood samples that are not similar, resulting in wrong pseudo-labels. On the contrary, when 481 non-parametric geometric feature is added to construct the neighborhood, the performance is usually 482 higher, and due to the stable geometric features on unseen data, the performance is still good even if 483 the number of neighborhoods is increased, which indicates the stability of our soft voting strategy. 484

Entropy Threshold. In this part, we explore the impact of σ and γ on performance which is 485 shown in Table 7 and Table 8. We use PointMetaBase-L as the source model and test it on the

- 473 474
- 475



Figure 3: Visualization of Feature Distribution on the 3DFRONT \rightarrow S3DIS benchmark. Front is the distribution of the source model with different methods. **Back** is the distribution of the model trained by labels on target domain.



Figure 4: Visualization of Segmentation Results on the 3DFRONT -> S3DIS benchmark.

3DFRONT \rightarrow S3DIS. In Table 7, we fixed γ =0.3 and observed that a smaller threshold effectively filters out difficult samples that are prone to generating erroneous predictions, thereby improving the model's performance. Similarly, in Table 8, we fixed σ = 0.2. And when set γ = 0, it indicates that all categories share the same threshold. We observed a significant drop in performance, which demonstrates the severe impact of class imbalance on performance.

Table	Table 7: Ablation Study of the σ						Table	e 8: Abla	ation Stu	dy of γ	
0.1	0.2	0.3	0.5	0.6			0	0.1	0.2	0.3	0.4
σ 36.18	36.36	36.01	35.91	35.52		$\gamma \mid$	34.92	35.63	35.85	36.36	36.12

4.5 VISUALIZATION

Feature Distribution. We visualized the distribution of features in Figure 3. In Figure 3(a), we visualize the distribution of source domains and target domains and see that they have a very large domain gap. In Figure 3(c), we visualize the non-parametric geometric feature distribution and the target domain distribution. We can see that they capture the information of the target domain distribution to a certain extent, such as the height between samples in their domain are relatively close. In Figure 3(d), the distribution of the source domain adjusted by our method is basically similar to that of the target domain, which effectively proves the effectiveness of our method.

530
 531
 531
 531
 531
 531
 531
 531
 531
 531
 531
 531
 531
 531
 531
 531
 531
 531
 531
 531
 531
 531
 531
 531
 531
 531
 531
 531
 531
 531
 531
 531
 531
 531
 531
 531
 531
 531
 531
 531
 531
 531
 531
 531
 531
 532
 532
 532
 532
 532
 532
 532
 532
 532
 532
 532
 532
 532
 532
 532
 532
 532
 532
 532
 532
 532
 532
 532
 532
 532
 532
 532
 532
 532
 532
 532
 532
 532
 532
 532
 532
 532
 532
 532
 532
 532
 532
 532
 532
 532
 532
 532
 533
 533
 534
 534
 534
 534
 534
 534
 534
 534
 534
 534
 534
 534
 534
 534
 534
 534
 534
 534
 534
 534
 534
 534
 534
 534
 534
 534

5 CONCLUSION

We argue that non-parameter geometric features can capture the underlying manifold of unseen data,
which has strong generalization. Therefore, we leverage non-parametric geometry as an intermediate
domain to prompt test time adaption. By introducing distillation from non-parametric model, pseudolabel refined by soft voting and category-balance sampler, our method can effectively improve the
performance of the source domain model in the target domain.

540 REFERENCES 541

582

542 543	Iro Armeni, Ozan Sener, Amir R Zamir, Helen Jiang, Ioannis Brilakis, Martin Fischer, and Silvio Savarasa. 3d semantic parsing of large scale indeer spaces. In <i>Proceedings of the IEEE Conference</i>
544	on Computer Vision and Pattern Recognition, pp. 1534–1543, 2016.
545 546 547	Dian Chen, Dequan Wang, Trevor Darrell, and Sayna Ebrahimi. Contrastive test-time adaptation. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp.

- 547 295-305, 2022. 548
- 549 Shuhao Cui, Shuhui Wang, Junbao Zhuo, Chi Su, Qingming Huang, and Qi Tian. Gradually 550 vanishing bridge for adversarial domain adaptation. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 12455–12464, 2020. 551

- 552 Angela Dai, Angel X. Chang, Manolis Savva, Maciej Halber, Thomas Funkhouser, and Matthias 553 Nießner. Scannet: Richly-annotated 3d reconstructions of indoor scenes. In Proc. Computer Vision 554 and Pattern Recognition (CVPR), IEEE, 2017. 555
- Runyu Ding, Jihan Yang, Li Jiang, and Xiaojuan Qi. Doda: Data-oriented sim-to-real domain 556 adaptation for 3d semantic segmentation. In European Conference on Computer Vision, pp. 284-303. Springer, 2022. 558
- 559 Huan Fu, Bowen Cai, Lin Gao, Ling-Xiao Zhang, Jiaming Wang, Cao Li, Qixun Zeng, Chengyue 560 Sun, Rongfei Jia, Binqiang Zhao, et al. 3d-front: 3d furnished rooms with layouts and semantics. 561 In Proceedings of the IEEE/CVF International Conference on Computer Vision, pp. 10933–10942, 2021a. 562
- 563 Huan Fu, Rongfei Jia, Lin Gao, Mingming Gong, Binqiang Zhao, Steve Maybank, and Dacheng Tao. 3d-future: 3d furniture shape with texture. International Journal of Computer Vision, 129: 565 3313-3337, 2021b. 566
- Robert Geirhos, Carlos RM Temme, Jonas Rauber, Heiko H Schütt, Matthias Bethge, and Felix A 567 Wichmann. Generalisation in humans and deep neural networks. Advances in Neural Information 568 Processing Systems, 31, 2018. 569
- 570 Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image 571 recognition. In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pp. 770-778, 2016. 572
- 573 Dan Hendrycks and Thomas Dietterich. Benchmarking neural network robustness to common corrup-574 tions and perturbations. Proceedings of the International Conference on Learning Representations, 575 2019. 576
- Ji Hou, Benjamin Graham, Matthias Nießner, and Saining Xie. Exploring data-efficient 3d scene 577 understanding with contrastive scene contexts. In Proceedings of the IEEE/CVF conference on 578 computer vision and pattern recognition, pp. 15587–15597, 2021. 579
- 580 Yusuke Iwasawa and Yutaka Matsuo. Test-time classifier adjustment module for model-agnostic 581 domain generalization. Advances in Neural Information Processing Systems, 34:2427–2440, 2021.
- Maximilian Jaritz, Tuan-Hung Vu, Raoul de Charette, Emilie Wirbel, and Patrick Pérez. xmuda: 583 Cross-modal unsupervised domain adaptation for 3d semantic segmentation. In Proceedings of the 584 IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 12605–12614, 2020. 585
- 586 Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolu-587 tional neural networks. Communications of the ACM, 60(6):84–90, 2017.
- 588 Shuang Li, Mixue Xie, Kaixiong Gong, Chi Harold Liu, Yulin Wang, and Wei Li. Transferable 589 semantic augmentation for domain adaptation. In Proceedings of the IEEE/CVF Conference on 590 Computer Vision and Pattern Recognition, pp. 11516–11525, 2021. 591
- Haojia Lin, Xiawu Zheng, Lijiang Li, Fei Chao, Shanshan Wang, Yan Wang, Yonghong Tian, 592 and Rongrong Ji. Meta architecture for point cloud analysis. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 17682–17691, 2023.

594 595 596	Mingsheng Long, Han Zhu, Jianmin Wang, and Michael I Jordan. Deep transfer learning with joint adaptation networks. In <i>International Conference on Machine Learning</i> , pp. 2208–2217. PMLR, 2017.
597	
598	Shuaicheng Niu, Jiaxiang Wu, Yifan Zhang, Yaofo Chen, Shijian Zheng, Peilin Zhao, and Mingkui
599	Tan. Efficient test-time model adaptation without forgetting. In International Conference on
600	Machine Learning, pp. 16888–16905. PMLR, 2022.
601	Shuaichang Niu, Jiaviang Wu, Vifan Zhang, Zhiguan Wan, Vaofo Chan, Dailin Zhao, and Mingkui
602	Tan Towards stable test-time adaptation in dynamic wild world. In <i>Internetional Conference on</i>
603	Learning Representations, 2023.
604	
605	Duo Peng, Yinjie Lei, Wen Li, Pingping Zhang, and Yulan Guo. Sparse-to-dense feature matching:
606	Intra and inter domain cross-modal learning in domain adaptation for 3d semantic segmentation.
607 608	In Proceedings of the IEEE/CVF International Conference on Computer Vision, pp. 7108–7117, 2021.
609	
610	Charles Ruizhongtai Qi, Li Yi, Hao Su, and Leonidas J Guibas. Pointnet++: Deep hierarchical feature
611	tearning on point sets in a metric space. Advances in Neural Information Processing Systems, 30, 2017
612	2017.
613	Guocheng Qian, Yuchen Li, Houwen Peng, Jinjie Mai, Hasan Hammoud, Mohamed Elhoseiny, and
614	Bernard Ghanem. Pointnext: Revisiting pointnet++ with improved training and scaling strategies.
615	Advances in Neural Information Processing Systems, 35:23192–23204, 2022.
616	
617	Haoxi Ran, Jun Liu, and Chengjie Wang. Surface representation for point clouds. In <i>Proceedings of</i>
618	the IEEE/CVF Conference on Computer vision and Pattern Recognition, pp. 18942–18952, 2022.
619	Benjamin Recht, Rebecca Roelofs, Ludwig Schmidt, and Vaishaal Shankar. Do imagenet classifiers
620	generalize to imagenet? In International Conference on Machine Learning, pp. 5389–5400. PMLR,
621	2019.
622	
623	Kuniaki Saito, Yoshitaka Ushiku, and Tatsuya Harada. Asymmetric tri-training for unsupervised
624	aomain adaptation. In International Conference on Machine Learning, pp. 2988–2997. PMLR, 2017
625	2017.
626	Inkyu Shin, Yi-Hsuan Tsai, Bingbing Zhuang, Samuel Schulter, Buyu Liu, Sparsh Garg, In So Kweon,
627	and Kuk-Jin Yoon. Mm-tta: Multi-modal test-time adaptation for 3d semantic segmentation. In
628	Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp.
629	16928–16937, 2022.
630	Junha Song Kwanyong Park InKyu Shin Sanghyun Woo Chaoning Zhang and In So Kweon
031	Test-time adaptation in the dynamic world with compound domain knowledge management. <i>IEEE</i>
622	Robotics and Automation Letters, 2023.
62/	
625	Shuofeng Sun, Yongming Rao, Jiwen Lu, and Haibin Yan. X-3d: Explicit 3d structure modeling for
635	point cloud recognition. In Proceedings of the IEEE/CVF Conference on Computer Vision and
637	Pattern Recognition, pp. 50/4-5085, 2024.
638	Hugues Thomas, Charles R Oi, Jean-Emmanuel Deschaud, Beatriz Marcotegui, Francois Goulette
639	and Leonidas J Guibas. Kpconv: Flexible and deformable convolution for point clouds. In
640	Proceedings of the IEEE/CVF International Conference on Computer Vision, pp. 6411–6420,
641	2019.
642	
643	Y1-HSuan 1sai, Wei-Unih Hung, Samuel Schulter, Kihyuk Sohn, Ming-Hsuan Yang, and Manmohan Chandraker, Learning to adapt structured output grace for computing segmentation. In Proceedings
644	of the IFFF Conference on Computer Vision and Pattern Recognition on 7472–7481 2018
645	of the HEEE conference on computer vision and ration Recognition, pp. 1412-1401, 2010.
646	Dequan Wang, Evan Shelhamer, Shaoteng Liu, Bruno Olshausen, and Trevor Darrell. Tent: Fully test-
647	time adaptation by entropy minimization. In <i>International Conference on Learning Representations</i> , 2021. URL https://openreview.net/forum?id=uXl3bZLkr3c.

- 648 Qin Wang, Olga Fink, Luc Van Gool, and Dengxin Dai. Continual test-time domain adaptation. 649 In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 650 7201-7211, 2022. 651
- 652 Zicheng Wang, Zhen Zhao, Yiming Wu, Luping Zhou, and Dong Xu. Progressive target-styled 653 feature augmentation for unsupervised domain adaptation on point clouds. arXiv preprint arXiv:2311.16474, 2023. 654
- Bichen Wu, Xuanyu Zhou, Sicheng Zhao, Xiangyu Yue, and Kurt Keutzer. Squeezesegv2: Improved 656 model structure and unsupervised domain adaptation for road-object segmentation from a lidar 657 point cloud. In 2019 International Conference on Robotics and Automation, pp. 4376–4382. IEEE, 658 2019. 659
 - Xiaoyang Wu, Xin Wen, Xihui Liu, and Hengshuang Zhao. Masked scene contrast: A scalable framework for unsupervised 3d representation learning. In Proceedings of the IEEE/CVF Conference on computer vision and pattern recognition, pp. 9415–9424, 2023.
 - Jihan Yang, Ruijia Xu, Ruiyu Li, Xiaojuan Qi, Xiaoyong Shen, Guanbin Li, and Liang Lin. An adversarial perturbation oriented domain adaptation approach for semantic segmentation. In Proceedings of the AAAI Conference on Artificial Intelligence, volume 34, pp. 12613–12620, 2020.
- Marvin Zhang, Sergey Levine, and Chelsea Finn. Memo: Test time robustness via adaptation and 668 augmentation. Advances in Neural Information Processing Systems, 35:38629–38642, 2022a. 669
- Renrui Zhang, Ziyu Guo, Peng Gao, Rongyao Fang, Bin Zhao, Dong Wang, Yu Qiao, and Hong-671 sheng Li. Point-m2ae: multi-scale masked autoencoders for hierarchical point cloud pre-training. 672 Advances in neural information processing systems, 35:27061–27074, 2022b. 673
- 674 Renrui Zhang, Liuhui Wang, Yali Wang, Peng Gao, Hongsheng Li, and Jianbo Shi. Starting from 675 non-parametric networks for 3d point cloud analysis. In Proceedings of the IEEE/CVF Conference 676 on Computer Vision and Pattern Recognition, pp. 5344–5353, 2023. 677
- 678 Bowen Zhao, Chen Chen, and Shu-Tao Xia. DELTA: DEGRADATION-FREE FULLY TEST-TIME 679 ADAPTATION. In The Eleventh International Conference on Learning Representations, 2023. 680 URL https://openreview.net/forum?id=eGm22rqG93.
- 682 Hengshuang Zhao, Li Jiang, Jiaya Jia, Philip HS Torr, and Vladlen Koltun. Point transformer. In 683 Proceedings of the IEEE/CVF International Conference on Computer Vision, pp. 16259–16268, 2021.
 - Yang Zou, Zhiding Yu, BVK Kumar, and Jinsong Wang. Unsupervised domain adaptation for semantic segmentation via class-balanced self-training. In Proceedings of the European Conference on Computer Vision (ECCV), pp. 289-305, 2018.
- 688 689 690

691 692

694

681

684 685

686

687

655

660

661

662

663

665

666

667

670

А APPENDIX

693 A.1 LIGHTWEIGHT POINTNN

To reduce computational overhead, we modified the settings of PointNN by decreasing the number 695 of layers (L), the number of neighboring points (K), downsample ratio (\mathbf{R}), and the init feature 696 dimensions C to enhance the efficiency of NGTTA. Although the modifications made to PointNN 697 may result in a slight decrease in accuracy, the resulting improvements in efficiency are significant. 698 This enhances the applicability of NGTTA in real-world scenarios. 699

- 700 701
- PointNN: L=5,K=90,R=2,C=144
- Ours: L=4,K=24,R=4,C=36

A.2 ABLATION STUDY OF THE LOSS WEIGHT α

The overall loss definition of NGTTA is as follows:

$$L_{tta} = \alpha_1 L_{dis} + \alpha_2 L_{cls} + \alpha_3 L_{ent} \tag{1}$$

where α_1, α_2 and α_3 mean the loss weight of L_{dis}, L_{cls} and L_{ent} .

709 In the following three tables, we use PointMetaBase-L as the source model and test the impact of the 10ss weight on the 3DFRONT→ S3DIS. From the results, we can see that when the weight is set to 0, 11there is a decline in performance, which demonstrates the necessity of each loss component. At the 11there is a decline in performance does not fluctuate significantly around the optimal weight 11there is a decline that our method is not particularly sensitive to the loss weights, thereby proving the 11there is a baseline that our method.

Table 1	: A	blat	tion S	Study	of α_1	Tab	le 2: A	blation	Study	of α_2	Table	e 3: Ał	olation	Study	of α_3
	0)	1	10	100		0	0.1	0.5	1		0	0.1	0.5	1
$alpha_1$	33.	81	34.96	36.36	36.01	α_2	34.10	35.89	36.36	36.21	α_3	34.60	35.91	36.30	36.36

A.3 ABLATION STUDY OF OPTIMIZABLE PARAMETERS

722 In our implementation, we only optimized the initial layers of the encoder and the classification layer. 723 In this section, we explore the impact of different optimization parameters on performance. We set up 724 various combinations of optimizable parameters: 1) Optimize only batch normalization parameters 725 (BN) 2) Optimize the entire model (AM) 3) Optimize only the classifier (CH) 4) Optimize only the 726 encoder (EC) 5) Optimize only the decoder (DC) 6) Our implementation (Ours) 7) Optimize the 727 encoder and the decoder (BC). We use PointMetaBase-L as the source model and conduct experiment 728 on 3DFRONT \rightarrow S3DIS. From the results, we can see that due to the high complexity of the semantic segmentation model, updating too many parameters can actually lead to a decline in performance, 729 as observed in the cases of AM and BC. In contrast, updating only a small number of important 730 parameters, such as in our method or BN, can achieve better results. 731

 Table 4: Ablation Study of the optimizable parameters

	BN	AM	СН	EC	DC	Ours	BC
mIOU	35.72	34.29	35.62	35.52	34.71	36.36	34.78

732 733

734 735 736

705 706

715

720

721

A.4 CATEGORY-WISE RESULTS

741 742 A.4.1 CATEGORY-WISE ENTROPY.

The class imbalance issue is very serious in point cloud scene segmentation, leading to significant differences in average entropy across different classes. To address this, we propose the Category-Balance Sampler. Here, we report the average entropy for each class for 3DFRONT \rightarrow S3DIS and 3DFRONT \rightarrow ScanNet, with the results shown in Table 5. It can be observed that minority classes typically have higher average entropy, which necessitates a higher entropy threshold to select more samples. This supports the validity of our Category-Balance Sampler.

749 750

751

A.4.2 CATEGORY-WISE IOU.

To further explore how our method improves the source domain model, we report the IoU for
 3DFRONT→ ScanNet, with the results shown in Table 6. It can be observed that a significant
 improvement of NGTTA lies in its ability to enhance the performance of minority and difficult classes
 effectively. In contrast, TENT's inability to address the class imbalance issue results in a decline in
 the performance of minority classes, leading to training collapse.

Table 5: Category-Wise Entropy Result on ScanNet and S3DIS.

						G					
					(a) ScanNet					
	wall	floor	cabinet	bed	chair	sofa	table	door	window	bookshelf	desk
Quantity Ratio	20%	20%	3%	2%	6%	2%	3%	3%	2%	2%	1%
Entropy	0.40	0.50	0.70	0.66	0.57	0.56	0.58	0.53	0.93	0.88	0.98
					(t	o) S3DIS					
	wall	floor	chair	sofa	table	door	window	bookshelf	ceiling	beam	column
Quantity Ratio	27%	15%	2%	1%	3%	3%	3%	11%	19%	1%	1%
Entropy	0.14	0.19	0.71	0.41	0.89	0.59	0.50	0.22	0.15	0.60	0.58

Table 6: Category-Wise IoU Result on 3DFRONT -> ScanNet.

Method wall	floor	cabinet	bed	chair	sofa	table	door	window	bookshelf	desk mIoU
Baseline 60.80	83.22	13.15	47.93	56.35	47.38	39.61	1.85	3.28	18.95	21.07 34.80
TENT 47.12	60.55	0.03	12.79	4.35	2.29	30.22	0	0	0	0 15.63
NGTTA 60.12	86.35	9.29	42.75	57.97	44.87	44.85	3.20	7.13	31.26	27.04 38.42

A.5 ADDITIONAL EXPERIMENTS

776 A.5.1 SEMANTICKITTI.

Here, we introduce a more challenging experiment by transferring from indoor data (3DFRONT-SemanticKITTI) to outdoor data to demonstrate the generalization capability of NGTTA. However, a significant challenge arises because indoor and outdoor datasets do not share the same classes, which prevents the use of common technical components such as pseudo-labeling, entropy minimization, and subsequent evaluation phases. Therefore, we only utilize a non-parametric geometric model for feature distillation, extracting point-level features from both the source model and the adaptive model, and we use SVM to evaluate accuracy. The results are shown in Table 7, where we can see that the accuracy significantly improves after using NGTTA. This demonstrates that NGTTA is also applicable to outdoor data and can enhance feature distinguishability.

Table 7: SVM Accuracy of NGTTA on 3DFRONT-SemanticKITTI.

	PointMetaBase-L	+NGTTA
Accuracy	30.5	35.1

A.5.2 COMPARE TO FPFH.

Here, we compare the performance of the geometric feature FPFH and PointNN on $3DFRONT \rightarrow ScanNet$ and $3DFRONT \rightarrow S3DIS$ to demonstrate that, in our method, PointNN is a superior non-parametric geometric feature extractor. The results are shown in Table 8. It can be observed that PointNN outperforms FPFH. We believe this is primarily because FPFH calculates geometric features based solely on the local neighborhood of each point, resulting in a very limited receptive field. In contrast, PointNN expands the receptive field by aggregating geometric features through multiple layers of down-sampling and up-sampling, which is crucial for scene segmentation.

803 804	Table 8: Results of FPFH.		
805		3DFRONT	3DFRONT→ScanNet
807	PointMetaBase-L	29.38	34.80
808	+FPFH	34.98	36.76
809	+PointNN	36.36	38.42

A.5.3 RESULTS ON ADDITIONAL MODELS.

in this part, we have added experiments with additional models MinkowskiNet (ResNet-UNet) and
RandLA-Net on 3DFRONT→S3DIS and 3DFRONT→ScanNet, and the results are shown in Table
As can be seen, NGTTA can be applied to various point cloud architectures, demonstrating the
generalizability of our method.

Table 9: Results of MinkowskiNet and RandLA-Net.

	3DFRONT->S3DIS	3DFRONT→ScanNet
MinkowskiNet	24.35	31.72
+NGTTA	32.57	35.76
RandLA-Net	25.81	32.16
+NGTTA	33.71	35.91

A.6 VISUALIZATION

827 A.6.1 SEGMENTATION RESULTS

In this part, we compare the visualization results of the Source Model and the results after applying
 NGTTA adaptation, as shown in Figure 5. It can be seen that NGTTA effectively corrects the
 erroneous class predictions of the Source Model, resulting in improved segmentation results.

A.6.2 ADAPTION PROCESS

Here, we visualize the domain adaptation process using NGTTA, with the results shown in Figure 6.
(a) represents the feature distribution obtained by directly applying the source model on the target data, while (d) represents the feature distribution of the model trained on labeled data in the target domain, which can be considered as the target domain distribution. (b) and (c) show the results after training with NGTTA for 1 and 2 epochs, respectively. We can see that the black and red circles highlight the areas where the source domain distribution and the target domain distribution differ significantly. As the NGTTA training progresses, the model's feature distribution gradually approaches the target domain distribution.

A.7 INTRODUCTION OF POINTNN

PointNN is a non-parametric geometric model that utilizes common components from point cloud models, such as Farthest Point Sampling (FPS), kNN, and max pooling, to extract local features from point clouds. Specifically, for the *i*-th point $p_i = (x_i, y_i, z_i) \in \mathbb{R}^{1 \times 3}$ in the point cloud, it first employs trigonometric functions to extract positional features.

851

852

853

854 855 856

859 860 861

862

842 843

844

834

816 817

$$f_i^x = [sine(Ax_i/B^{\frac{6\cdot 0}{C}}), Cosine(Ax_i/B^{\frac{6\cdot 0}{C}}), ..., sine(Ax_i/B^{\frac{6\cdot m}{C}}), Cosine(Ax_i/B^{\frac{6\cdot m}{C}})]$$
(2)

where f_i^x is the positional encoding for the x-axis. A and B means the magnitude and wavelengths. The encoding for the y-axis and z-axis is the same. Therefore, the position embedding of p_i can be written as:

$$f_{i} = PoE(p_{i}) = [f_{i}^{x}, f_{i}^{y}, f_{i}^{z}]$$
(3)

Then, PointNN use kNN to construct the neighborhood $(p_j, f_j)_{j \in \mathcal{N}_i}$ of the *i*-th center point. Then, for each neighborhood vector, PointNN expands it to:

$$f_{ij} = [f_i, f_j] \tag{4}$$

To capture the relevant geometric information between the center point and the neighboring points, PointNN incorporates relative positional encoding:



Figure 5: Visualization comparison of NGTTA and Source Model. The first column is the label, the second column is the segmentation result from the Source Model, and the last column is the segmentation result after applying NGTTA adaptation. Then, the first two rows are the results for $3DFRONT \rightarrow ScanNet$, and the last two rows are the results for $3DFRONT \rightarrow S2DIS$.

$$\hat{f}_{ij} = (f_{ij} + PoE((p_i - p_j))) \odot PoE((p_i - p_j))$$
(5)

$$\hat{f}_i = MaxPool(\{\hat{f}_{ij}\}_{j \in \mathcal{N}_i}) + AvgPool(\{\hat{f}_{ij}\}_{j \in \mathcal{N}_i})$$
(6)

PointNN extracts rich geometric information from local regions by employing multi-layer downsampling and aggregating local features and then obtains point-wise features through multi-layer upsampling.



Figure 6: TSNE of NGTTA Results on 3DFRONT-S3DIS. The black and red circles represent the parts where the source model distribution and the target domain distribution differ significantly. It can be observed that as NGTTA continues to train, the model's feature distribution gradually approaches the target domain distribution.