

AutoGPart: Intermediate Supervision Search for Generalizable 3D Part Segmentation

Xueyi Liu¹ Xiaomeng Xu¹ Anyi Rao² Chuang Gan³ Li Yi^{1,4}
¹ Tsinghua University ² The Chinese University of Hong Kong ³ MIT-IBM Watson AI Lab
⁴ Shanghai Qi Zhi Institute

Abstract

Training a generalizable 3D part segmentation network is quite challenging but of great importance in real-world applications. To tackle this problem, some works design task-specific solutions by translating human understanding of the task to machine’s learning process, which faces the risk of missing the optimal strategy since machines do not necessarily understand in the exact human way. Others try to use conventional task-agnostic approaches designed for domain generalization problems with no task prior knowledge considered. To solve the above issues, we propose AutoGPart, a generic method enabling training generalizable 3D part segmentation networks with the task prior considered. AutoGPart builds a supervision space with geometric prior knowledge encoded, and lets the machine to search for the optimal supervisions from the space for a specific segmentation task automatically. Extensive experiments on three generalizable 3D part segmentation tasks are conducted to demonstrate the effectiveness and versatility of AutoGPart. We demonstrate that the performance of segmentation networks using simple backbones can be significantly improved when trained with supervisions searched by our method. Project page: <https://autogpart.github.io>.

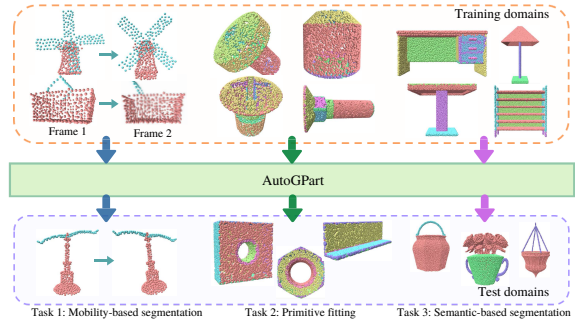


Figure 1. AutoGPart automatically finds real part cues for different segmentation tasks and supports simple point cloud processing backbones (e.g. PointNet++ [36]) to successfully parse shapes from novel categories.

same distribution as the training set. However, they usually suffer from a huge performance drop when it comes to parsing 3D shapes from a novel distribution, e.g. from a different semantic category.

In this work, we focus on generalizable 3D part segmentation tasks. The goal is to learn the essence of parts on some training sets and to be able to generalize well to shapes from novel distributions. A main challenge comes from the versatile cues defining parts. Shape parts can be defined via many cues like geometric primitive fitting [15, 24, 47], rigid motions [13, 14, 48], and semantic prior [28, 30, 33]. The same shape can be segmented into different parts based upon different cues for the seek of various applications. Therefore for a specific application, other than understanding input shapes, a network also needs to figure out the exact cues defining parts from the training data. This is where it becomes tricky. Some of these cues are more salient than others, making them easier to be captured. When these cues heavily correlate with the real part cue in the training domain, they could easily become shortcut features [8, 9, 16] biasing the network and hindering the generalizability, as observed in [30].

To cope with the above challenges, existing works usually focus on a specific type of parts and incorporate part-type priors towards generalizable 3D part segmentation. For example, a modularized design has been adopted to explicitly extract motion flow and the rigid

1. Introduction

Humans parse objects into parts for a deeper understanding of semantics, functionality, mobility as well as the fabrication process. It is natural to consider equipping machines with such part-level understanding. Over the past few years, there has been an increasing interest in 3D part segmentation tasks to support various applications [15, 24, 33, 47] in vision, graphics, and robotics. Thanks to the availability of large-scale 3D part datasets [33, 49] and the development of 3D deep learning techniques, these methods achieve impressive part segmentation results when a test shape comes from the

motion of local patches so that generalizable motion-based part segmentation can be achieved in [48]. Similarly, the networks are guided explicitly to extract primitive parameters in order to segment primitive parts in [24, 47]. Finding such explicit part type priors requires a deep understanding of the target task, which might not always be available. Even after a huge amount of trial-and-errors in expert designs, such approaches are still not guaranteed to fully grasp the essence of parts. This is because human understandings might not align well with the way that machines prefer to follow [7, 19, 53].

Another line of works from the machine learning community studies generic domain generalization algorithms to help a network handle out-of-distribution samples. However, without considering any geometric prior which is important for part segmentation tasks, they fail to improve crucial parts of the model but merely focus on generic regularization strategies such as making gradients from multiple domains consistent with each other [31].

In contrast, we present *AutoGPart*, a generalizable 3D part segmentation technique that could be applied to any type of 3D parts. Our key observation is that the generalizability of a 3D part segmentation network is largely hindered by shortcut features [8]. These shortcut features tend to be quite salient and closely correlated with segmentation labels in some training domains. Through finding proper intermediate supervisions that are more closely related with real part cues and loosely correlated with shortcut features, we can downweight the influence of shortcut features by jointly training the network with segmentation and the intermediate supervisions.

We design a parametric model to depict the distribution of possible intermediate supervisions. And we inject geometric priors in the model space such that these supervisions are part-aware and geometrically-discriminative. In addition, we present an automatic way for optimizing the distribution of these supervisions for a supervised part segmentation network. A beam search-like strategy is then applied to greedily generate some suitable supervisions from the distribution. With the additional supervision, task-specific part cues can be automatically discovered without experts' trail-and-errors. Comprehensive experiments on three generalizable 3D part segmentation tasks demonstrate the effectiveness of *AutoGPart*, including 4.4% absolute Segmentation Mean IoU improvement on mobility-based part segmentation, 4.2% on primitive fitting and 0.7% absolute Mean Recall [30] improvement on semantic-based part segmentation.

To summarize, our contributions are threefold: 1) We present a generic method to improve the generalizability of 3D part segmentation networks via automatically discovered intermediate supervisions and the method is

suitable for different part definitions. 2) We design a parametric model to depict the distribution of useful intermediate supervisions with geometric prior encoded. 3) We propose an automatic search algorithm to find the proper intermediate supervision given a specific part segmentation task.

2. Related Work

Domain Generalization. As an important technique to handle out-of-distribution scenarios, existing domain generalization approaches can be categorized into four streams: 1) learning domain invariant features from multiple source domains aiming to minimize the difference between source domains [23, 26, 34], meta-learning [21, 22], and other model-agnostic strategies [31]; 2) data augmentation algorithms to simulate domain shift [56]; 3) ensemble learning techniques that train domain-specific models and uses their ensemble for inference [46, 56]; 4) automated machine learning (AutoML) techniques that aim to automatically search for data augmentation [5, 25] or neural architecture [2, 3] which can achieve optimal generalization performance. Different from previous AutoML strategies, we propose to search for geometric and part-aware features for intermediate supervisions. Such features are invariant across shapes from different distributions, indicating the superiority of our method to improve the generalizability of the network.

3D Part Segmentation. Mobility-based part segmentation, primitive fitting, and semantic-based part segmentation are three important and representative 3D part segmentation tasks, on which we conduct experiments to prove the effectiveness and versatility of *AutoGPart*. 1) *Mobility-based part segmentation* aims to parse articulated input objects into rigidly moving parts. A number of previous works have devoted into it based on traditional techniques such as clustering and co-segmentation [42, 51], or deep 3D neural networks [14, 48]. 2) *Primitive fitting* addresses the task of clustering input points and fitting them with geometric primitives. Standard solutions include RANSAC [38], region growing [32], supervised learning [15, 24, 47] and unsupervised learning [6, 41]. Recently, [24] proposes an end-to-end neural network that takes point clouds as input and predicts a varying number of primitives. [47] further uses hybrid feature representations to separate points of different primitives. 3) *Semantic-based part segmentation* detects and delineates each distinct object of interest. Conventional approaches rely on manual design on geometric constraints, including K-means [39], graph cuts [10] and spectral clustering [27]. As for learning-based methods [35, 50], although they achieve state-of-the-art performance on seen categories, it is hard for them to parse shapes from unseen categories due to category-variant information that is easy to take

than real cues defining parts [30]. Different from previous task-specific methods that employ sophisticated frameworks and delicately designed supervisions, we demonstrate that simple learning frameworks equipped with the intermediate supervisions searched from AutoGPart can achieve comparable or even better performance than existing task-specific designs on the three tasks.

Auxiliary Supervision. Adding auxiliary supervisions is a common strategy to regularize the learning process for performance improvement. Among them, weight decay is a widely used technique [18, 29]. Besides, multi-task learning [4, 12, 37, 52] trains a network to solve multiple tasks by sharing parameters between different tasks, where the interested branches acquire benefits from other branches. Deeply Supervise [20] and Inception [40] explore how to add auxiliary supervisions on hidden layers to improve the quality of learned low-level features [54, 55]; LabelEnc [11] proposes a new label encoding function that maps ground-truth labels to the latent embedding space to add intermediate supervisions. In this work, we automatically select proper intermediate supervisions for generalizable 3D part segmentation networks. It can be viewed as an approach that utilizes auxiliary supervisions to improve the network’s generalizability.

3. Method

We present AutoGPart, a technique to improve the cross-domain generalizability of a supervised 3D part segmentation network by automatically finding useful intermediate supervisions for the network. Then at training time, supervisions found by AutoGPart are introduced in addition to the segmentation supervision to alleviate the influence of shortcut features [8] and capture real part cues that are invariant across different shapes distributions. The resulting network can better generalize to shapes from novel distributions without any additional cost at inference time. Figure 2 shows an example where AutoGPart is plugged into an end-to-end network for primitive segmentation.

To automatically find such supervisions, we design a parametric supervision feature space with geometric prior knowledge encoded (*a parametric intermediate supervision model*, Sec. 3.1). To adapt the model for a specific segmentation network, we optimize its parameters via a “propose, evaluate, update” strategy (Sec. 3.2). After that, a greedy search strategy is utilized to select the optimal supervisions from the optimized model (Sec. 3.3).

3.1. Parametric Modeling for the Supervision Space

Based on the observation that many previous generalizable 3D part segmentation works guide part feature learning via intermediate supervisions calculated from per-point geometric features and ground-truth part

labels, we assume that being aware of some of such features can help the network learn real part cues. We propose to automatically search for proper part-aware supervision features for each type of part and add intermediate supervisions to encourage the segmentation network to leverage those features. We then move on to introduce the structural model for supervision features and further the supervision space.

Structural model for supervision features. We define the possible supervision features as the outcome of a tree-structured operation flow (an *operation tree*) taking geometric features and ground-truth part labels as input. It is inspired by the calculation process of hand-crafted part-aware geometric features (*e.g.* rotation matrix, cylinder axis [24], etc.).

Good supervision features could bridge the gap between input point cloud geometry and the output segmentation labels while avoiding shortcuts. And it is crucial to consider how to use ground-truth labels in order to generate such **part-aware and geometrically-discriminative** supervision features. To compute a possible intermediate supervision, instead of treating ground-truth labels as additional input feature vectors, we pre-encode these labels by transferring geometric features of each point to part-aware features. We then use such part-aware features as input to an operation tree directly.

An operation tree transforms input features by operators for the output intermediate supervision feature. Such operators include *grouping operators* (*e.g.* sum, SVD¹, etc.) that summarize a feature set into a point-level feature, point-level *unary operators* (*e.g.* square, double, etc.), and *binary operators* (*e.g.* add, minus, etc.) that encourage feature communications to enlarge and diversify the supervision space. We further introduce some fixed operator-type combinations named *operation cells* such as a unary operator followed by a grouping operator. They are high-frequent operation combinations in the computing process of hand-crafted geometric features [24, 48].

An operation tree is then constructed by connecting operation cells.

Supervision feature space. The supervision feature space consists of all valid operation trees. We set the maximum height of an operation tree to three, thus the supervision feature space is spanned by all possible tree structures and sub-structures of cells in it (Figure 3). To measure the generalization benefit of each supervision feature and to optimize the space toward the optimal intermediate supervisions for a segmentation network, a parametric distribution model ($\mathcal{M}_{\mathcal{T}}(\cdot|\theta)$) is constructed to depict the operation tree space. The subscription \mathcal{T} here indicates

¹Singular Value Decompose.

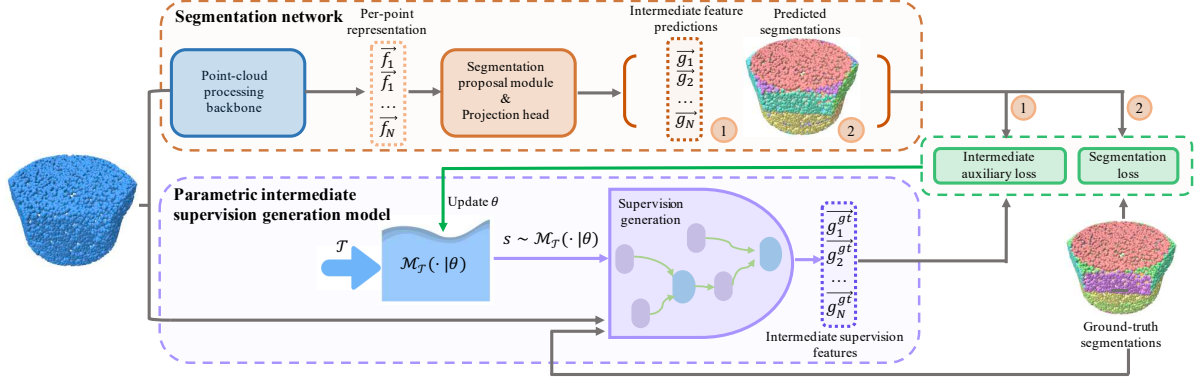


Figure 2. An example of applying AutoGPart on an end-to-end segmentation network for the primitive fitting task. AutoGPart builds a learnable parametric supervision model $\mathcal{M}_{\mathcal{T}}(\cdot|\theta)$ that can be optimized in order to find proper intermediate supervisions for a generalizable 3D part segmentation network. A “propose, evaluate and update” strategy is adopted to learn θ . In each circle, an operation tree is sampled from \mathcal{M} , which is then used to calculate part-aware geometric features for each point. The generalization ability of the proposed supervision is evaluated and the resulting score is further used to update θ in the updating stage. Light purple lines imply the supervision generation process and green for the parameter updating process.

that it is constructed based on prior knowledge of the task set \mathcal{T} , θ denotes parameters introduced in the distribution model. In practice, a marginal distribution is introduced for the grouping operator of the root cell. Then, conditional distributions are introduced for children operators and leaf part-aware features conditioned on their parents.

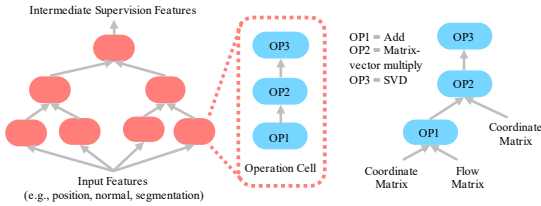


Figure 3. **Left side:** Left side: Supervision feature space. **Right Side:** An example for the operation tree whose output feature is the rotation matrix calculated by the Orthogonal Procrustes algorithm. Operation cell abstraction is not drawn on it.

3.2. Learning Task-Conditioned Supervision Distribution

To optimize the distribution parameters θ for a segmentation network, we adopt a “propose, evaluate and update” strategy.

Supervision proposal. Since a supervision feature is modeled by its generation process in our model, an operation tree depicting the generation of a supervision feature is sampled from the parametric model $\mathcal{M}_{\mathcal{T}}(\cdot|\theta)$ to generate a supervision feature. After that, the input part-aware geometric feature matrices are passed through the sampled operation tree to generate the corresponding supervision feature s . Then the ground-truth intermediate supervision feature \vec{g}_i^{gt} for each point i is calculated by passing its input features through s .

Supervision evaluation. The cross-domain generalization ability of the proposed supervision feature is estimated under a simulated domain-shift setting. Firstly, the training domain is split into different sub-domains. After that, the OOD performance of the sampled supervision feature is estimated by the average generalization error calculated on all of those train-validation splits created by learning one sub-domain out for validation. The network is trained together with the proposed supervision s and the segmentation task-related supervision to evaluate the generalization benefit of s .

Supervision space optimization. The parameter θ is updated by the probability density value of the generated supervision ($\mathcal{M}_{\mathcal{T}}(s|\theta)$) and its estimated generalization score. The updating strategy is derived from the REINFORCE [45] algorithm. Simply updating parameters of conditional distributions involved in the generation process of s via REINFORCE equals to updating the joint distribution by REINFORCE.

3.3. Greedy Supervision Selection

Though we can select a single supervision feature from the supervision model after the optimization (e.g. the feature with the highest probability), we wish to select multiple features to use for better generalizability enhancement.

To be more specific, we adopt a greedy search strategy to select several supervisions (one to three, in our practice) to use further for the segmentation network. The greedy search aims to choose an optimal supervision set step-by-step. It starts from a set containing single supervisions sampled from the optimized supervision space. The generalization ability of such supervision features are then estimated. After that, supervisions with

top performance are kept and further used to construct a set containing two supervision pairs. Similarly, the performance of supervision combinations are estimated with a new set constructed by coupling top supervisions in each following step. Finally, the supervision combination achieved the best performance is then selected to use further.

4. Experiments

We evaluate the effectiveness of the proposed AutoGPart in searching for suitable intermediate supervisions and improving the domain generalization ability of a segmentation network on three 3D part segmentation tasks: mobility-based part segmentation [14, 48], primitive fitting [15, 24, 47] and semantic-based part segmentation [30].

For each task, we set a default segmentation network to evaluate AutoGPart on, which is a point-cloud processing backbone for representation learning followed by a classification-based segmentation module to propose segmentations. Two backbones are used in experiments, namely PointNet++ [36] and DGCNN [44]. We denote such two segmentation networks as “PointNet++” and “DGCNN” directly for simplicity. For each task, we demonstrate that a simple 3D point-cloud segmentation network trained with intermediate supervisions searched by AutoGPart and segmentation task-related supervisions is able to perform better than task-specific models proposed in previous works as well as networks trained by generic domain generalization strategies.

4.1. Mobility-based Part Segmentation

Datasets. Three datasets are used in the task: training dataset, auxiliary training dataset that is only used in the supervision search stage to help simulate domain-shift, and the out-of-domain (OOD) test dataset. The training dataset is created from [49], containing 15,776 shapes from 16 categories. The auxiliary training dataset is created from PartNet [33], containing 5,662 shapes from 4 categories different from those used in the training dataset. The test dataset used is the same as the one used in [48], which is created from [14], containing 875 shapes covering 175 objects from 23 categories.

Experimental settings. We evaluate the effectiveness of AutoGPart on both PointNet++ and DGCNN. Metric used for this task is Segmentation Mean IoU (MIoU). To apply AutoGPart on this task, we add a flow estimation module ahead of the segmentation network. We compare our method with both task-specific baselines [42, 48, 51], and task-agnostic methods [21, 56]. Deep Part Induction [48] is trained on the same training dataset using the same settings as those used for our models. All learning-based models

Table 1. Experimental results on the mobility-based segmentation task. For abbreviations used, In/Out-of-dist. refers to In/Out-of-distribution performance; “PN++” denotes “PointNet++”; “JLC” means “JLinkage clustering”; “SC” means “Spectral Clustering”. Subscripts indicate the used backbones.

Method	In-dist.	Out-of-dist.
PointNet++	86.4	61.0
DGCNN [44]	88.3	68.1
MixStyle _{PN++} [56]	86.6	63.4
MixStyle _{DGCNN} [56]	83.1	69.5
Meta-learning _{PN++} [21]	71.3	63.4
Meta-learning _{DGCNN} [21]	76.5	69.4
JLC [51]	N/A	67.3
SC [42]	N/A	69.4
Deep Part Induction [48]	84.8	64.4
AutoGPart _{PN++}	87.2	66.5
AutoGPart _{DGCNN}	83.1	73.8

use only one forward flow estimation and segmentation proposal pass with no iteration between them.

Experimental results. The performance of AutoGPart and baseline models are summarized in Table 1. Based on the table, we can make the following observations: 1) AutoGPart can find useful intermediate supervisions that can significantly improve the generalization ability of both PointNet++ and DGCNN (*e.g.* 4.9% absolute MIoU improvement on OOD-test set for PointNet++). This can verify the effectiveness of adding intermediate supervisions to boost the generalization ability of segmentation networks and the ability of AutoGPart to find such supervisions. 2) AutoGPart can outperform all task-specific models by a large margin, including traditional methods (*e.g.* JLinkage clustering) and learning based strategies such as Deep Part Induction. This may echo the proposed assumption that hand-crafted supervisions may not be optimal ones for the task due to the misalignment between human knowledge and machine understanding. 3) AutoGPart can outperform all task-agnostic strategies when using the same point-cloud processing backbone. A possible reason is that geometric prior knowledge which is quite useful to solve the task inherently is carefully considered in the design of AutoGPart, while task-agnostic strategies are not aware of such prior knowledge.

4.2. Primitive Fitting

Dataset. Dataset used in this task is the same as the one used in [24]. However, instead of using the provided data splitting method directly, we re-split the dataset into 4 subsets according to shapes such that it is more suitable to test the cross-domain generalization ability of a model. Three subsets out of them, containing 13,528 shapes in total, are used for training, including the supervision search stage and regular training stage for supervision evaluation.

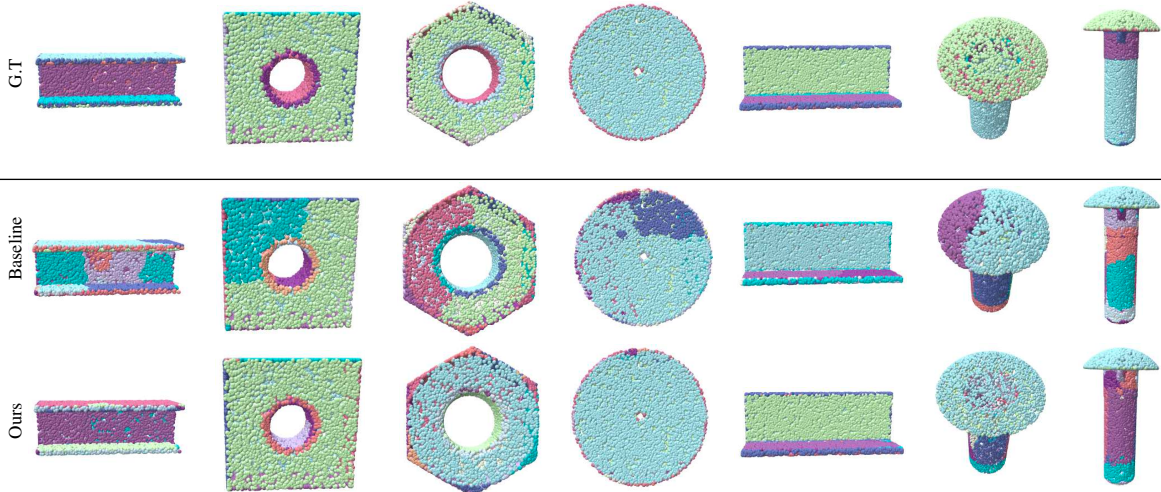


Figure 4. Segmentation results visualization for the primitive fitting task. “Ours” denotes the model “AutoGPart_{HPNet}” and “Baseline” denotes the model “HPNet” in Table 2.

The left one subset, containing 3,669 shapes, is used for the out-of-domain test. For all compared baselines, we test their performance on the re-split dataset for a fair comparison.

Experimental settings. We evaluate the effectiveness of AutoGPart on both PointNet++ and DGCNN. Metric used for this task is MIOU. We compare our method with both task-specific baselines [24, 47] and task-agnostic methods [21, 31, 56].

Among them, HPNet [47] adopts a clustering-based two-stage segmentation network for feature learning and segmentation proposal. Based on the observation that the high performance achieved by HPNet is largely powered by its clustering-based segmentation module and the synergies between the added supervisions and the clustering module, we conduct the following two experiments to fairly compare AutoGPart with HPNet, since the classification-based segmentation module used in default settings of AutoGPart is not that representative compared with Mean-Shift clustering used in HPNet: 1) Plug AutoGPart in HPNet’s first learning stage to evaluate the effectiveness of AutoGPart in finding useful features that can regularize this learning stage. 2) Replace the clustering-based segmentation proposal module used in HPNet with a classification-based module to compare the effectiveness of supervisions introduced in HPNet and those searched by AutoGPart under a same segmentation network. Resulted models are denoted as “AutoGPart_{HPNet}” and “HPNet*” respectively in Table 2.

Experimental results. The performance of AutoGPart and other baseline models are summarized in Table 2. We can make the following observations from Table 2: 1) AutoGPart can find useful intermediate supervisions to improve the generalization ability of a simple segmentation network for the primitive fitting task using either

PointNet++ or DGCNN as the backbone, similar with that observed for mobility-based part segmentation task. 2) For models using classification-based segmentation modules, AutoGPart can boost a simple segmentation network’s performance better than all task-specific models such as SPFN. A possible reason is that the added supervisions, though believed useful by human to help the network learn a correct solution inherently, may not align well with the way preferred by machines to solve the task. Thus, those hand-crafted supervisions may not be effective enough to train a generalizable network. Besides, networks may even use shortcut features to optimize those supervisions. For models using clustering-based segmentation modules, AutoGPart_{HPNet} can achieve better generalization performance than HPNet by adding some intermediate supervisions in its first learning stage. It is probably because that the added supervisions searched by AutoGPart can help the network use more real cues to learn per-point features in this stage. 3) AutoGPart can help better improve the generalization ability of 3D part segmentation networks compared with task-agnostic methods.

The absolute improvement that AutoGPart adds on HPNet is not as significant as PointNet++ or DGCNN, for which we want to emphasize two points as follows: 1) The high performance achieved by HPNet is largely benefit from the clustering-based segmentation module (*i.e.* 69.6% MIOU on OOD-test set achieved by HPNet* *v.s.* 79.5% achieved by HPNet). However, such a segmentation proposal process is quite **time-consuming**, where about 40 hours are needed to segment all test shapes, while a classification-based one only needs no more than 40 seconds. 2) For classification-based segmentation networks, intermediate supervisions searched by AutoGPart are clearly better than that used in HPNet

Table 2. Experimental results on the primitive fitting task. For abbreviations used, In/Out-of-dist. refers to In/Out-of-distribution performance; “PN++” denotes “PointNet++”. Subscriptions indicate the used backbones; “GS” means “Gradient Surgery”. “HPNet” and “AutoGPart_{HPNet}” use clustering-based segmentation modules. Others use classification-based modules.

Method	In-dist.	Out-of-dist.
PointNet++	81.5	71.6
DGCNN	93.6	68.0
GS _{PN++} [31]	90.7	70.3
GS _{DGCNN} [31]	92.6	71.6
Meta Learning _{PN++} [21]	65.0	68.5
Meta-learning _{DGCNN} [21]	67.1	69.3
MixStyle _{PN++} [56]	92.1	70.9
MixStyle _{DGCNN} [56]	93.7	71.4
SPFN [24]	94.4	72.3
HPNet* [47]	93.9	69.6
HPNet [47]	N/A	79.5
AutoGPart _{PN++}	86.3	76.5
AutoGPart _{DGCNN}	94.2	73.4
AutoGPart _{HPNet}	N/A	80.4

by comparing the performance of AutoGPart_{DGCNN} with HPNet* (*i.e.* 73.4% *v.s.* 69.6% MIoU on OOD-test set). The significant effectiveness of AutoGPart in improving the generalization ability of an end-to-end trained model is of larger practical value.

4.3. Semantic-based Part Segmentation

Dataset. We use the same dataset provided by [30] as well as the train-test data splitting strategy stated in the paper.

Experimental Settings. The evaluation metric used in this task is the Mean Recall value, the same as the one used in [30]. Values reported in Table 3 are the average Mean Recall scores over all test categories. We compare AutoGPart with four task-specific methods [17, 30, 43, 50], and two task-agnostic approaches [21, 56].

Experimental results. The performance of AutoGPart and other baseline models are summarized in Table 3. Different from the above two tasks, where the criterion about segmenting a shape is obvious such as rigid motions, it is not that clear what part cues are useful for this task, meaning not enough prior knowledge to guide a learning-based network’s design. Thus, the role of ground-truth geometric features is not that important in the previous task-specific designs [30, 43, 50]. However, their resulting models tend to perform not that well when parsing shapes from novel categories, as shown in Table 3, probably due to the influence of shortcut features. However, our method can help improve the generalization ability of a simple segmentation network to a level comparable to the two-stage learning-based method [30] as well as traditional segmentation methods

Table 3. Experimental results on the semantic-based part segmentation task. For abbreviations used, In/Out-of-dist. refers to In/Out-of-distribution performance. The backbone used in AutoGPart is PointNet++.

Method	In-dist.	Out-of-dist.
MixStyle [56]	34.4	30.7
Meta-learning [21]	35.1	29.7
PartNet-InstSeg [33]	33.9	26.7
SGPN [43]	25.0	20.2
GSPN [50]	23.5	28.7
Learning to Group [30]	35.2	32.0
WCseg [17]	29.8	33.2
AutoGPart	35.7	33.9

such as WCseg [17]. It demonstrates the usefulness of ground-truth part-aware features in providing real cues for a generalizable segmentation network; And also indicates the superiority of AutoGPart to find such useful features for a task where human prior knowledge is not that available or hard to be translated to guide a network’s learning process.

4.4. Qualitative Evaluation

We visualize the segmentation results of “AutoGPart_{HPNet}” and “HPNet” in Table 2 on the Primitive Fitting task for a intuitive comparison and understanding w.r.t. the network’s generalization ability on OOD shapes. Figure 4 shows that AutoGPart_{HPNet} can achieve better segmentation performance on shapes from a novel distribution, *i.e.*, having a relatively large primitive-type distribution gap from that of training shapes. A possible reason is that the added intermediate supervisions in the first stage of HPNet can help the model learn per-point features relying more on real part cues, thus less sensitive to the primitive type of each part as well as the primitive-type distribution of the shape.

Besides, we draw an intermediate supervision feature searched from our designed supervision space in Figure 6 for an intuitive understanding towards the property of our searched feature. It can be seen that the searched feature is not only discriminative across different parts, but also changes continuously within one part. Such property may make it be more friendly for the network learning real part cues avoiding shortcut features that only exist in shapes from training distributions.

5. Ablation Study

In our method, we construct an intermediate supervision space and search for useful supervisions from it to increase the generalizability of a segmentation network. In this section, we ablate our method by replacing some crucial designs with other possible alternatives to analyze the effect of those parts.

Table 4. Ablation study w.r.t. reward function design and supervision space design. For abbreviations used, “Arch.” refers to “Architecture”; “PN++” denotes “PointNet++”; In/Out-of-dist. refers to In/Out-of-distribution performance.

Ablation	Arch.	In-dist.	Out-of-dist.
/	PN++	87.2	66.5
–Cross Validation		85.6	64.6
–Gap		90.6	64.9
Less operands		87.5	63.1
Less unary operators		86.8	64.1
Less binary operators		83.0	65.1
Less tree height		82.1	63.5
/	DGCNN	83.1	73.8
–Cross Validation		84.5	73.1
–Gap		89.4	70.4
Less operands		92.1	70.6
Less unary operators		89.3	69.5
Less binary operators		89.7	71.6
Less tree height		89.4	71.8

Reward function design. In AutoGPart, we adopt the average generalization gap over all train-validation splits as the reward value to estimate the effectiveness of the selected supervision on improving the model’s generalization ability. The intuition is that the cross-validating the generalization gap of the searched supervision feature could probably give its generalizability a better estimation. In Table 4, we validate the superiority of the designed reward function by replacing it with 1) generalization gap on a single split (denoted as “–Cross validation”), and 2) average performance on validation sets over all splits (denoted as “–Gap”).

Supervision space design. In our method, we use a large supervision feature space to increase the diversity of features in it and also to enable the model a high freedom to select its preferred features by itself at the same time. To demonstrate its benefit, we downsize the supervision space by downsizing the input part-aware geometric feature set, the unary operator set and the binary operator set respectively and then evaluate the effectiveness of the searched supervisions. Experimental results prove the superiority of using a large supervision space (see Table 4).

Supervision selection strategy. In our method, we use a greedy search-like strategy to select supervision features for future use based on the intuition that such a strategy can help us find a good supervision feature combination from the optimized space. To demonstrate the effectiveness of the greedy supervision selection strategy used in our method, we try to ablate it and consider the following two alternatives: 1) Select top K ($1 \leq K \leq 3$) features from the supervision feature set ordered by the probability density value. Results are summarized in Table 5. 2) Choose features from the optimized distributions randomly. Ten random features are sampled and their results are summarized in Figure 5.

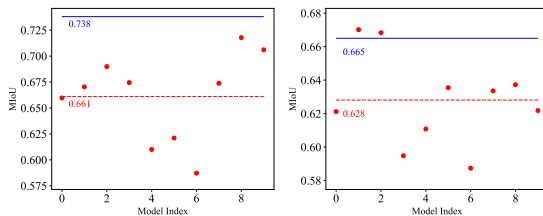


Figure 5. Ablation study w.r.t. supervision selection strategy. Ten supervision features are randomly selected from the optimized distributions and evaluated. Left for PointNet++ and right for DGCNN.

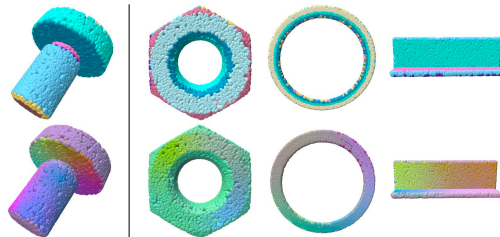


Figure 6. Visualization for one searched intermediate supervision feature from our designed feature space. Upper: Segmentations; Lower: Searched feature (converted to RGB values). The left most shape is the one from the training dataset, while the other three are from the test dataset.

Table 5. Ablation study w.r.t. the supervision selection strategy. For abbreviations used, “Arch.” refers to “Architecture”; “PN++” denotes “PointNet++”; In/Out-of-dist. refers to In/Out-of-distribution performance.

Ablation	Arch.	In-dist.	Out-of-dist.
/	PN++	87.2	66.5
Top1		87.1	65.6
Top2		81.5	64.5
Top3		85.8	62.8
/	DGCNN	83.1	73.8
Top1		92.1	70.1
Top2		89.8	67.8
Top3		90.7	71.6

6. Discussion and Conclusion

In this paper, we propose to automatically search for proper intermediate supervisions for generalizable 3D part segmentation networks.

Although experimental results can prove the effectiveness and versatility of AutoGPart to some extent, there are still many directions worth exploring by the community in the future: 1) Search intermediate supervisions and the backbone’s structure at the same time. Synergies between backbone architecture and intermediate supervisions may be found since different backbones may prefer different features. 2) How to add intermediate supervisions. Directly predicting a ground-truth feature may not be suitable for each feature. Moreover, synergies between supervisions features and how to supervise the network to learn such features may be discovered.

References

- [1] David Arthur and Sergei Vassilvitskii. k-means++: The advantages of careful seeding. Technical report, Stanford, 2006. [15](#)
- [2] Haoyue Bai, Fengwei Zhou, Lanqing Hong, Nanyang Ye, S-H Gary Chan, and Zhenguo Li. Nas-ood: Neural architecture search for out-of-distribution generalization. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 8320–8329, 2021. [2](#)
- [3] Corinna Cortes, Xavier Gonzalvo, Vitaly Kuznetsov, Mehryar Mohri, and Scott Yang. AdaNet: Adaptive structural learning of artificial neural networks. In Doina Precup and Yee Whye Teh, editors, *Proceedings of the 34th International Conference on Machine Learning*, volume 70 of *Proceedings of Machine Learning Research*, pages 874–883. PMLR, 06–11 Aug 2017. [2](#)
- [4] Michael Crawshaw. Multi-task learning with deep neural networks: A survey. *arXiv preprint arXiv:2009.09796*, 2020. [3](#)
- [5] Ekin D Cubuk, Barret Zoph, Dandelion Mane, Vijay Vasudevan, and Quoc V Le. Autoaugment: Learning augmentation strategies from data. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 113–123, 2019. [2](#)
- [6] Hao Fang, Florent Lafarge, and Mathieu Desbrun. Planar shape detection at structural scales. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2965–2973, 2018. [2](#)
- [7] Christina M Funke, Judy Borowski, Karolina Stosio, Wieland Brendel, Thomas SA Wallis, and Matthias Bethge. The notorious difficulty of comparing human and machine perception. *arXiv e-prints*, pages arXiv–2004, 2020. [2](#)
- [8] Robert Geirhos, Jörn-Henrik Jacobsen, Claudio Michaelis, Richard Zemel, Wieland Brendel, Matthias Bethge, and Felix A Wichmann. Shortcut learning in deep neural networks. *Nature Machine Intelligence*, 2(11):665–673, 2020. [1](#), [2](#), [3](#), [16](#)
- [9] Robert Geirhos, Patricia Rubisch, Claudio Michaelis, Matthias Bethge, Felix A Wichmann, and Wieland Brendel. Imagenet-trained cnns are biased towards texture; increasing shape bias improves accuracy and robustness. *arXiv preprint arXiv:1811.12231*, 2018. [1](#)
- [10] Aleksey Golovinskiy and Thomas Funkhouser. Randomized cuts for 3d mesh analysis. In *ACM SIGGRAPH Asia 2008 papers*, pages 1–12. 2008. [2](#)
- [11] Miao Hao, Yitao Liu, Xiangyu Zhang, and Jian Sun. Labelenc: A new intermediate supervision method for object detection. In *European Conference on Computer Vision*, pages 529–545. Springer, 2020. [3](#)
- [12] Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross Girshick. Mask r-cnn. In *Proceedings of the IEEE international conference on computer vision*, pages 2961–2969, 2017. [3](#)
- [13] Ruizhen Hu, Lubin Fan, and Ligang Liu. Co-segmentation of 3d shapes via subspace clustering. In *Computer graphics forum*, volume 31, pages 1703–1713. Wiley Online Library, 2012. [1](#)
- [14] Ruizhen Hu, Wenchao Li, Oliver Van Kaick, Ariel Shamir, Hao Zhang, and Hui Huang. Learning to predict part mobility from a single static snapshot. *ACM Transactions on Graphics (TOG)*, 36(6):1–13, 2017. [1](#), [2](#), [5](#), [15](#)
- [15] Jingwei Huang, Yanfeng Zhang, and Mingwei Sun. Primitivenet: Primitive instance segmentation with local primitive embedding under adversarial metric. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 15343–15353, 2021. [1](#), [2](#), [5](#)
- [16] Andrew Ilyas, Shibani Santurkar, Dimitris Tsipras, Logan Engstrom, Brandon Tran, and Aleksander Madry. Adversarial examples are not bugs, they are features. *arXiv preprint arXiv:1905.02175*, 2019. [1](#)
- [17] Oliver Van Kaick, Noa Fish, Yanir Kleiman, Shmuel Asafi, and Daniel Cohen-Or. Shape segmentation by approximate convexity analysis. *ACM Transactions on Graphics (TOG)*, 34(1):1–11, 2014. [7](#), [17](#)
- [18] Anders Krogh and John A Hertz. A simple weight decay can improve generalization. In *Advances in neural information processing systems*, pages 950–957, 1992. [3](#)
- [19] Sebastian Lapuschkin, Stephan Wäldchen, Alexander Binder, Grégoire Montavon, Wojciech Samek, and Klaus-Robert Müller. Unmasking clever hans predictors and assessing what machines really learn. *Nature communications*, 10(1):1–8, 2019. [2](#)
- [20] Chen-Yu Lee, Saining Xie, Patrick Gallagher, Zhengyou Zhang, and Zhuowen Tu. Deeply-supervised nets. In *Artificial intelligence and statistics*, pages 562–570. PMLR, 2015. [3](#)
- [21] Da Li, Yongxin Yang, Yi-Zhe Song, and Timothy M Hospedales. Learning to generalize: Meta-learning for domain generalization. In *Thirty-Second AAAI Conference on Artificial Intelligence*, 2018. [2](#), [5](#), [6](#), [7](#), [16](#)
- [22] Da Li, Jianshu Zhang, Yongxin Yang, Cong Liu, Yi-Zhe Song, and Timothy M Hospedales. Episodic training for domain generalization. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 1446–1455, 2019. [2](#)
- [23] Haoliang Li, Sinno Jialin Pan, Shiqi Wang, and Alex C Kot. Domain generalization with adversarial feature learning. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 5400–5409, 2018. [2](#)
- [24] Lingxiao Li, Minhyuk Sung, Anastasia Dubrovina, Li Yi, and Leonidas J Guibas. Supervised fitting of geometric primitives to 3d point clouds. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2652–2660, 2019. [1](#), [2](#), [3](#), [5](#), [6](#), [7](#), [15](#), [16](#)
- [25] Yonggang Li, Guosheng Hu, Yongtao Wang, Timothy Hospedales, Neil M Robertson, and Yongxin Yang. Dada: Differentiable automatic data augmentation. *arXiv preprint arXiv:2003.03780*, 2020. [2](#)
- [26] Ya Li, Xinmei Tian, Mingming Gong, Yajing Liu, Tongliang Liu, Kun Zhang, and Dacheng Tao. Deep domain generalization via conditional invariant adversarial networks. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 624–639, 2018. [2](#)

- [27] Rong Liu and Hao Zhang. Segmentation of 3d meshes through spectral clustering. In *12th Pacific Conference on Computer Graphics and Applications, 2004. PG 2004. Proceedings.*, pages 298–305. IEEE, 2004. 2
- [28] Shu Liu, Jiaya Jia, Sanja Fidler, and Raquel Urtasun. Sgn: Sequential grouping networks for instance segmentation. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 3496–3504, 2017. 1
- [29] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. *arXiv preprint arXiv:1711.05101*, 2017. 3
- [30] Tiange Luo, Kaichun Mo, Zhiao Huang, Jiarui Xu, Siyu Hu, Liwei Wang, and Hao Su. Learning to group: A bottom-up framework for 3d part discovery in unseen categories. *arXiv preprint arXiv:2002.06478*, 2020. 1, 2, 3, 5, 7, 15, 17
- [31] Lucas Mansilla, Rodrigo Echeveste, Diego H. Milone, and Enzo Ferrante. Domain generalization via gradient surgery. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 6630–6638, October 2021. 2, 6, 7, 16
- [32] David Marshall, Gabor Lukacs, and Ralph Martin. Robust segmentation of primitives from range data in the presence of geometric degeneracy. *IEEE Transactions on pattern analysis and machine intelligence*, 23(3):304–314, 2001. 2
- [33] Kaichun Mo, Shilin Zhu, Angel X Chang, Li Yi, Subarna Tripathi, Leonidas J Guibas, and Hao Su. Partnet: A large-scale benchmark for fine-grained and hierarchical part-level 3d object understanding. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 909–918, 2019. 1, 5, 7, 15
- [34] Krikamol Muandet, David Balduzzi, and Bernhard Schölkopf. Domain generalization via invariant feature representation. In Sanjoy Dasgupta and David McAllester, editors, *Proceedings of the 30th International Conference on Machine Learning*, volume 28 of *Proceedings of Machine Learning Research*, pages 10–18, Atlanta, Georgia, USA, 17–19 Jun 2013. PMLR. 2
- [35] Charles R Qi, Hao Su, Kaichun Mo, and Leonidas J Guibas. Pointnet: Deep learning on point sets for 3d classification and segmentation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 652–660, 2017. 2
- [36] Charles R Qi, Li Yi, Hao Su, and Leonidas J Guibas. Pointnet++: Deep hierarchical feature learning on point sets in a metric space. *arXiv preprint arXiv:1706.02413*, 2017. 1, 5, 17
- [37] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. Faster r-cnn: Towards real-time object detection with region proposal networks. *Advances in neural information processing systems*, 28:91–99, 2015. 3
- [38] Ruwen Schnabel, Roland Wahl, and Reinhard Klein. Efficient ransac for point-cloud shape detection. In *Computer graphics forum*, volume 26, pages 214–226. Wiley Online Library, 2007. 2
- [39] Shymon Shlafman, Ayellet Tal, and Sagi Katz. Metamorphosis of polyhedral surfaces using decomposition. In *Computer graphics forum*, volume 21, pages 219–228. Wiley Online Library, 2002. 2
- [40] Christian Szegedy, Wei Liu, Yangqing Jia, Pierre Sermanet, Scott Reed, Dragomir Anguelov, Dumitru Erhan, Vincent Vanhoucke, and Andrew Rabinovich. Going deeper with convolutions. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1–9, 2015. 3
- [41] Shubham Tulsiani, Hao Su, Leonidas J Guibas, Alexei A Efros, and Jitendra Malik. Learning shape abstractions by assembling volumetric primitives. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2635–2643, 2017. 2
- [42] Dimitrios Tzionas and Juergen Gall. Reconstructing articulated rigged models from rgb-d videos. In *European Conference on Computer Vision*, pages 620–633. Springer, 2016. 2, 5, 16
- [43] Weiyue Wang, Ronald Yu, Qiangui Huang, and Ulrich Neumann. Sgpn: Similarity group proposal network for 3d point cloud instance segmentation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2569–2578, 2018. 7, 17
- [44] Yue Wang, Yongbin Sun, Ziwei Liu, Sanjay E Sarma, Michael M Bronstein, and Justin M Solomon. Dynamic graph cnn for learning on point clouds. *Acm Transactions On Graphics (tog)*, 38(5):1–12, 2019. 5
- [45] Ronald J Williams. Simple statistical gradient-following algorithms for connectionist reinforcement learning. *Machine learning*, 8(3):229–256, 1992. 4
- [46] Zheng Xu, Wen Li, Li Niu, and Dong Xu. Exploiting low-rank structure from latent domains for domain generalization. In *European Conference on Computer Vision*, pages 628–643. Springer, 2014. 2
- [47] Siming Yan, Zhenpei Yang, Chongyang Ma, Haibin Huang, Etienne Vouga, and Qixing Huang. Hpnet: Deep primitive segmentation using hybrid representations. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 2753–2762, October 2021. 1, 2, 5, 6, 7, 16, 17
- [48] Li Yi, Haibin Huang, Difan Liu, Evangelos Kalogerakis, Hao Su, and Leonidas Guibas. Deep part induction from articulated object pairs. *arXiv preprint arXiv:1809.07417*, 2018. 1, 2, 3, 5, 15, 16
- [49] Li Yi, Vladimir G Kim, Duygu Ceylan, I-Chao Shen, Mengyan Yan, Hao Su, Cewu Lu, Qixing Huang, Alla Sheffer, and Leonidas Guibas. A scalable active framework for region annotation in 3d shape collections. *ACM Transactions on Graphics (ToG)*, 35(6):1–12, 2016. 1, 5, 15
- [50] Li Yi, Wang Zhao, He Wang, Minhyuk Sung, and Leonidas J Guibas. Gspn: Generative shape proposal network for 3d instance segmentation in point cloud. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3947–3956, 2019. 2, 7, 17
- [51] Qing Yuan, Guiqing Li, Kai Xu, Xudong Chen, and Hui Huang. Space-time co-segmentation of articulated point cloud sequences. In *Computer Graphics Forum*, volume 35, pages 419–429. Wiley Online Library, 2016. 2, 5, 16
- [52] Yu Zhang and Qiang Yang. A survey on multi-task learning. *arXiv preprint arXiv:1707.08114*, 2017. 3

- [53] Zijian Zhang, Jaspreet Singh, Ujwal Gadiraju, and Avishek Anand. Dissonance between human and machine understanding. *Proceedings of the ACM on Human-Computer Interaction*, 3(CSCW):1–23, 2019. [2](#)
- [54] Zhenli Zhang, Xiangyu Zhang, Chao Peng, Xiangyang Xue, and Jian Sun. Exfuse: Enhancing feature fusion for semantic segmentation. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 269–284, 2018. [3](#)
- [55] Hengshuang Zhao, Jianping Shi, Xiaojuan Qi, Xiaogang Wang, and Jiaya Jia. Pyramid scene parsing network. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2881–2890, 2017. [3](#)
- [56] Kaiyang Zhou, Yongxin Yang, Yu Qiao, and Tao Xiang. Domain generalization with mixstyle. In *International Conference on Learning Representations*, 2021. [2](#), [5](#), [6](#), [7](#), [16](#)
- [57] Qian-Yi Zhou, Jaesik Park, and Vladlen Koltun. Open3d: A modern library for 3d data processing. *arXiv preprint arXiv:1801.09847*, 2018. [16](#)

A. Overview

The appendix aims at providing more details on the method design (Sec. B), experimental settings (Sec. C), and more experimental results (Sec. D).

B. Method Details

In this section, we give some explanations for some details of the proposed method that are not stated in the main paper.

B.1. Supervision Feature Structure

This section gives more details about designed structural model for generating ground-truth part-aware geometric features w.r.t. further explanations for each component, operation cells and the rationality of the design.

Input features. Input features for each point are ordered part-aware feature sets formed by different kinds of geometric features from points having the same part labels with the target point. A maximum radius is set between the sampled points and the target point considering the limited receptive field of some point-cloud processing backbones such as DGCNN. This value is set to 0.5 in our implementation.

The full set of ordered input feature sets are formed by first sampling such part-aware feature sets from input geometric features such as coordinate vectors, then expanding the resulting sets by adding feature sets calculated by simple point-level operations among them to it. For example, suppose each point i has two different input geometric features: the coordinate vector \vec{p}_i and the normal vector \vec{n}_i . Their ordered respective part-aware feature sets for point i , constructed by sampling same kind of features from points having the same part labels with it, are denoted as \mathcal{P}_i and \mathcal{N}_i respectively. In the first step, we can get $\{\mathcal{P}_i, \mathcal{N}_i\}$ for point i . Then after expanding the set, we get $\{\mathcal{P}_i, \mathcal{N}_i, \mathcal{P}_i - \mathcal{N}_i, \dots\}$, where $\mathcal{P}_i - \mathcal{N}_i$ refers to a feature set formed by element-wise minus between \mathcal{P}_i and \mathcal{N}_i . An ordered part-aware geometric input feature set such as \mathcal{P}_i is referred as an ‘‘operant’’ in the operation tree. In practice, an ordered feature set is formed to a matrix with each of its line a feature vector of a point, thus P_i for \mathcal{P}_i and N_i for \mathcal{N}_i . If each point has two different geometric features, the full set of input part-aware geometric matrices contains 7 elements: $\{P_i, N_i, P_i \cdot N_i, P_i + N_i, P_i - N_i, N_i - P_i, \text{cross_product}(N_i, P_i)\}$. If each point has one geometric feature, the full set of input part-aware geometric matrices contains 4 elements: $\{P_i, 2P_i, P_i^2, -P_i\}$.

Operators. Operators are introduced to transform the input features step by step for the output intermediate supervision feature. Such operators include grouping operators, point-level unary operators, and point-level

Table 6. Candidates of each kind of feature transformation operator. ‘‘SVD’’ denotes ‘‘Singular Value Decompose’’, ‘‘Identity’’ refers to no operations, where the input ordered feature set (matrix) is not transformed by any unary operator.

Operators	Choices
Grouping operators	Sum, Average, Maximum, SVD
Binary operators	Add, Minus, Multiply, Cross Product, Cartesian Product, Matrix-vector Product
Unary operators	Identity, Square, Double, Negative, Orthogonalize, Inverse

binary operators, as summarized in Table 6. Then the grouping operator set, the unary operator set, and the binary operator set are further referred as \mathcal{G} , \mathcal{U} , and \mathcal{B} respectively.

Some operators may seem confusing, for which we give them some brief explanations as follows:

- **Orthogonalize:** Given a matrix M , first calculate its singular vectors and singular values by $U, S, V^T = \text{SVD}(M)$, then take the matrix-vector multiplication between two singular vector matrices for the result: $\text{Orthogonalize}(M) = UV^T$.
- **Cartesian Product:** Given two feature sets from one point, \mathcal{S}_1 and \mathcal{S}_2 , calculate the cartesian product between them by $\mathcal{S}_1 \times \mathcal{S}_2 = \{(s_i, s_j) | s_i \in \mathcal{S}_1, s_j \in \mathcal{S}_2\}$.

Operation cells. We introduce several different kinds of operation cells with fixed operator combinations to encourage some fixed operation sequences. Details of such operator combinations for operation cells in different levels are listed as follows:

- **Level-3 cells (Top level cells):** a grouping operator followed by a unary operator.
- **Level-2 cells:** a grouping operator followed by a unary operator.
- **Level-1 cells:** an operant followed by a unary operator.

Justification for the designed operation cells. In the design of operation cells, we set some fixed operator sequences for them such as a grouping operator followed by a unary operator. The design of the structure of operation cells is heuristic but also reasonable. For example, a unary operator following a grouping operator can further transform the grouped feature, extracting part-level information from it. Such calculation routines are common in the calculation process of some features such as the rotation matrix R . To be more specific, several steps in the calculation process of R contain a combination of ‘‘Centralize’’ and ‘‘Sum’’, where the former is a unary operator and the later is a grouping operator. This is only one possibility we explored in our practice, we believe there could be many other strategies to design operation cells (e.g. other computing routines) or just constructing the operation tree without introducing operation cells. How to design

Table 7. Examples of hand-crafted ground-truth features used in previous works. Matrices in the table are part-level matrices formed by corresponding features from points in a part. For notations used, F is the flow matrix, P is the coordinate matrix, N is the normal matrix. \bar{M} denotes the centralized matrix of M . Note that the geometric feature listed in the left column of the table may not be the exact feature output by the calculation process listed in its corresponding right cell, but a part of the output feature such as a row vector of the matrix.

Feature	Calculation process
Rotation matrix	$\text{orth}(\text{sum}(\text{cartesian}(P + F, P)))$
Cone apex	$\text{sum}(\text{cartesian}(N^{+T}, \text{rowsum}(N \cdot P)))$
Cylinder axis	$\text{svd}(\bar{N})$
Sphere center	$\text{sum}(\text{cartesian}((-2\bar{P})^{+T}, \text{rowsum}(P^2)))$
Sphere radius	$\text{sqrt}(\text{mean}(P - \vec{c})^2)$
Plane normal	$\text{svd}(\bar{P})$

suitable operation cells, or how to add proper constraints for operation combinations in the operation space is interesting and worth exploring.

Connections with hand-crafted geometric features. The constructed supervision feature space contains many hand-crafted geometric features used in previous works. Some examples are summarized in Table 7. The calculation process uses matrix forms of the ordered geometric feature sets.

B.2. Supervision Feature Distribution Space

This section aims for some further explanations of the constructed supervision feature distribution space w.r.t. how we decompose it into a tree-structured distribution space and the conditional sampling process performed on the distribution tree.

Decomposed tree-structured distribution set. Instead of using one single distribution, we decompose the total/joint distribution into a tree-structured distributions set to better model the generation process of the operation tree. Specifically, we use a distribution cell to depict the space of an operation cell, where the operator sequence in the operation cell is generated by a set of conditional operator distributions. For instance, a distribution cell for an operation cell with the ordered operator sequence [a grouping operator, a unary operator] is composed of a grouping operator distribution and $|\mathcal{G}|$ unary operator distributions. Each possible grouping operator has its own unary operator distribution. Thus, the distribution cell for this specific operation cell can be organized into a distribution tree with the grouping operator distribution as the top distribution and a unary operator distribution attached to each element in the grouping operator.

To model possible connections between cells from adjacent two levels, we introduce a connection distribution for each possible operator sequence of the upper-level cell to depict each possible connection and the operator that

should be used for connection. For example, a Level-3 cell can be connected with two Level-1 cells with each kind of binary operator, or a Level-2 cell and a Level-1 cell with each kind of binary operator, or a Level-1 cell and a Level-2 cell with each kind of binary operator, or a single Level-1 cell with no binary operator. Thus, the total number of possible connections is $3|\mathcal{B}| + 1$. A distribution containing elements of this number is introduced for each possible operator sequence of the upper-level cell.

Conditional operation tree sampling. After the distribution space has been constructed, a conditional sampling process is adopted to sample an operation tree from the space. Specifically, to sample an operation cell from its respective distribution cell, its top operator is sampled at first. Then, the distribution for the following level operator of this top operator is chosen to continue the sampling process. Specifically, the sampling process for sampling an operation sequence $[\text{op}_1, \text{op}_2, \dots, \text{op}_k]$ is conducted by: Sample op_k , conditioned on op_k and sample op_{k-1}, \dots , conditioned on $\text{op}_k, \text{op}_{k-1}, \dots, \text{op}_2$ and sample op_1 . The probability density value of the sequence is then calculated by:

$$p([\text{op}_1, \text{op}_2, \dots, \text{op}_k]) = \quad (1)$$

$$p(\text{op}_k)p(\text{op}_{k-1}|\text{op}_k)\dots p(\text{op}_1|[\text{op}_k, \dots, \text{op}_2]). \quad (2)$$

Similarly, to sample an operation tree, the top operation cell is first sampled. Then, a connection is sampled conditioned on the sampled operation cell structure. The sampled connection determines what distribution cells to continue the sampling process and what binary operator to use for connection if needed. Thus, the operation tree can be sampled by such a top-down manner.

B.3. Supervision Feature Distribution Learning

In this section, we talk about how to learn parameters for distributions in the supervision feature space, including the supervision feature sampling process, its cross-domain generalization ability evaluation and the distribution updating process.

Supervision feature sampling. A supervision feature is generated by first sampling an operation tree from the constructed distribution space and then pass input part-aware geometric features of each point through the operation tree to get its calculated part-aware geometric feature. The operation tree is sampled by the conditional sampling process from the distribution space as state above.

Cross-validation setting. Each selected supervision is put under a cross-validation process to estimate its cross-domain generalization ability. To cross-validate the generalization ability of the selected supervision, we first split the train dataset into K subsets with a relatively

large distribution shift across them. Then, the selected supervision is tested on each $K - 1 : 1$ train-validation fold. The average value between the metric on the validation set and the train set is taken as the estimated score for its generalization ability. In the supervision search process, the generalization score is further taken as the reward for the selected supervision and used for the following supervision distribution space update process. The cross-validation procedure for generalization gap evaluation is summarized in Algorithm 1.

Supervision distribution space optimization. We optimize distributions in each layer related with the sampling process of the selected supervision via the REINFORCE algorithm. Optimizing those conditional distributions via REINFORCE results optimizing the overall/joint distribution via REINFORCE. The optimization strategy is derived from the REINFORCE algorithm, where each parameter w_i that counts in the sampling process is updated by minusing its corresponding $\Delta w_i = \alpha_i(r - b_i)e_i$, where $e_i = \partial \ln g / \partial w_i$, g is the probability density function. Thus, e_i for parameter w_i of a distribution can be derived as:

$$e_i = \frac{\partial g}{\partial w_i} \quad (3)$$

$$= \frac{\partial \ln(g(\mathcal{H})g(v|\mathcal{H})g(\mathcal{C} \setminus \{v\}|\mathcal{H})g(\mathcal{L}|\mathcal{H}, \mathcal{C}))}{\partial w_i} \quad (4)$$

$$= \frac{\partial \ln(g(\mathcal{H}) + \partial \ln(g(v|\mathcal{H})) + \partial \ln(g(\mathcal{L}|\mathcal{H}, v)))}{\partial w_i} \quad (5)$$

$$= \frac{\partial \ln(g(v|\mathcal{H}))}{\partial w_i}, \quad (6)$$

where $g(\mathcal{H})$ is the probability density of values sampled in upper layers, $g(\mathcal{L}|\mathcal{H}, \mathcal{C})$ is the probability density of values sampled in lower layers conditioned on values sampled in current layers \mathcal{C} and those sampled in upper layers \mathcal{H} , which are not relevant with the sampling distribution in the current layer resulting value v . \mathcal{H} denotes the set of operators in upper layers of the operation tree, \mathcal{C} means the set of operators in the current layer of the operation tree, while \mathcal{L} refers to the set of operators in lower layers. The conditional distribution for each sample $g(v|\mathcal{H})$ is a multinomial distribution, in our implementation, which can be simply calculated by the class function “log_prob” of the corresponding PyTorch Distribution module.

B.4. Greedy Supervision Feature Selection

This section aims for some further explanations of the unstated details of the greedy feature selection process.

After we have got the optimized supervision distribution set, we greedily select suitable supervisions from the optimized space. At most three supervisions are selected

from the optimized space. To complete the supervision selection process, we first select K supervisions from the optimized space by randomly sampling from the optimized supervision space. Then the performance of such K supervisions are evaluated under a simulated domain-shift setting. From the ranked set, top 2 supervisions are selected to combine with top $K/2$ supervisions and form the second supervision set. Supervision combinations in the second supervision set are further evaluated and ranked. Then, top 3 supervision combinations containing 2 supervisions are selected to combine with top $K/3$ single supervisions from the first supervision set. The resulting combination set is further evaluated and ranked. The supervision combination that achieves the best estimated performance is selected for further evaluation. The performance of the selected supervision combination is regarded as the performance of the optimized supervision distribution space. The training stage of the supervision evaluation process is referred as the “regular training stage”. This greedy selection process is also summarized in Algorithm 3.

Algorithm 1 Cross_Val. “TrainValidation($\cdot, \cdot, \cdot, \cdot$)” takes a model, the train dataset, the validation dataset and the number of training epochs n as input, trains the model for n epochs and returns the gap between the performance of the model achieved on the validation set and the training set at the best validation epoch.

Input: The model \mathcal{M} ; Split train set $\mathcal{S}_{sp} = \{\mathcal{S}_1, \mathcal{S}_2, \dots, \mathcal{S}_k\}$; Intermediate supervision feature t ; Epochs for training n .

Output: Estimated generalization score which is actually the average generalization gap across all train-validation splits.

```

1: scores  $\leftarrow$  []
2: for  $i = 1$  to  $|\mathcal{S}_{sp}|$  do
3:    $\mathcal{S}_{tr} \leftarrow \mathcal{S}_{sp} \setminus \{\mathcal{S}_i\}$ 
4:    $\mathcal{S}_{val} \leftarrow \{\mathcal{S}_i\}$ 
5:    $s_i \leftarrow$  TrainValidation( $\mathcal{M}, \mathcal{S}_{tr}, \mathcal{S}_{val}, n$ )
6:   scores.append( $s_i$ )
7:  $\bar{s} \leftarrow$  mean(scores)
8: return  $\bar{s}$ 

```

Algorithm 2 SortByCrossVal.

Input: The model \mathcal{M} ; Split train set $\mathcal{S}_{sp} = \{\mathcal{S}_1, \mathcal{S}_2, \dots, \mathcal{S}_k\}$; A set of intermediate supervision features $\mathcal{T} = \{t_1, t_2, \dots, t_K\}$.

Output: Sorted intermediate supervision features \mathcal{T}' .

```

1: for  $i = 1$  to  $|\mathcal{T}|$  do
2:    $s \leftarrow$  Cross.Val( $\mathcal{M}, \mathcal{T}_1[i], \mathcal{S}_{sp}$ )
3:    $\mathcal{T}_1[i] \leftarrow (\mathcal{T}_1[i], s)$ 
4:  $\mathcal{T}' \leftarrow$  sorted( $\mathcal{T}$ )
5: return  $\mathcal{T}'$ 

```

Algorithm 3 GreedySupervisionSelection.

Input: The model \mathcal{M} ; Split train set $\mathcal{S}_{sp} = \{\mathcal{S}_1, \mathcal{S}_2, \dots, \mathcal{S}_k\}$; A set of sampled intermediate supervision features $\mathcal{T} = \{t_1, t_2, \dots, t_K\}$.

Output: Selected intermediate supervision feature set $\mathcal{T}_s = \{t_{s1}, \dots, t_{sk}\}, 1 \leq k \leq 3$.

- 1: $\mathcal{T}_1 \leftarrow \mathcal{T}$
 - 2: $\mathcal{T}_1 \leftarrow \text{SortByCrossVal}(\mathcal{M}, \mathcal{T}_1, \mathcal{S}_{sp})$
 - 3: $\mathcal{T}_2 \leftarrow \mathcal{T}_1[1 : 2] \times \mathcal{T}_1[1 : \lfloor |\mathcal{T}_1|/2 \rfloor]$
 - 4: $\mathcal{T}_2 \leftarrow \text{SortByCrossVal}(\mathcal{M}, \mathcal{T}_2, \mathcal{S}_{sp})$
 - 5: $\mathcal{T}_3 \leftarrow \mathcal{T}_2[1 : 3] \times \mathcal{T}_1[1 : \lfloor |\mathcal{T}_1|/3 \rfloor]$
 - 6: $\mathcal{T}_3 \leftarrow \text{SortByCrossVal}(\mathcal{M}, \mathcal{T}_3, \mathcal{S}_{sp})$
 - 7: $\mathcal{T}_s \leftarrow \text{sorted}(\{\mathcal{T}_1[1], \mathcal{T}_2[1], \mathcal{T}_3[1]\})[1]$
 - 8: **return** \mathcal{T}_s
-

C. Experimental Details

In this section, we provide some detailed information about experimental settings that are not stated in the main paper, including datasets used in each segmentation task, detailed settings for each stage, implementation for baseline models, etc.

C.1. Dataset

Mobility-based part segmentation. For the training dataset and auxiliary training dataset, we first infer part mobility information for parts in each shape. Then, during the training process, we generate shape pairs for training based on such part mobility information. For the training dataset, which is created from [49], containing 15,776 shapes from 16 categories, we first infer mobility information for parts in a shape heuristically based on their semantic labels such that the generated mobility information for the part can align better with real scenarios. Specifically, some heuristic constraints are added to infer mobility information for each part such as the chair back can rotate around the chair seat but not chair legs. For the auxiliary dataset created from PartNet [33], the mobility information for parts in a shape is also inferred heuristically where only some general motion rules are added, which means that parts of different semantic labels share the same motion rules. The test dataset used in our work is the same as the one used in [48], which is created from [14].

Primitive fitting. We use the same dataset as the one provided by [24], but adopt a different data splitting strategy such that it is more suitable to test a model’s cross-domain generalization ability. Specifically, we split shapes via their primitive-type distributions. Primitive-type distribution reveals the ratio of each primitive-type in the shape calculated based on number of points belonged to each type of primitive. We first cluster all shapes into 7 clusters by K-Means++ [1], then merge them into 4 subsets with a relatively large distribution gap across them. “Train_1”, “Train_2”, and “Train_3” are all used in the supervision search process and regular training process.

In the supervision search process, such three train splits serve for distribution-shift simulation to cross-validate the cross-domain generalization ability of the selected supervision. In the regular training process, shapes in those three splits are merged together and further split into train-validation datasets via a ratio 9:1. The validation set is used for model selection.

Semantic-based part segmentation. The dataset we use is the same as the one used in [30]. We only use the finest segmentation level for training and evaluation other than using all available levels as does in [30].

C.2. Experimental Settings

Supervision search. For mobility-based part segmentation, 100 automatic search epochs are conducted. For primitive fitting and semantic-based part segmentation, 30 automatic search epochs are conducted. For all three tasks, 4 supervisions sampled and evaluated in each epoch. The supervision distribution space is optimized in each epoch. AdaM optimizer is used for segmentation networks for all three tasks, with $\beta = (0.9, 0.999)$, $\epsilon = 10^{-8}$, and weight decay ratio set to 10^{-4} . Batch size is set to 36 for mobility-based part segmentation task, 2 for primitive fitting and semantic-based part segmentation using DGCNN, 8 those two tasks using PointNet++.

Regular training. We use AdaM optimizer for segmentation networks for all three tasks, with $\beta = (0.9, 0.999)$, $\epsilon = 10^{-8}$, and weight decay ratio set to 10^{-4} . For mobility-based part segmentation networks, 400 epochs are performed with the epoch that the model achieves the best validation performance is take for inference. For primitive fitting and semantic-based part segmentation tasks, 200 epochs are performed using the same model selection strategy as the one for the mobility-based part segmentation networks. When using clustering-based segmentation module, 100 training epochs are conducted with the model achieves the lowest validation contrastive-style loss further used for inference. As for the cross-validation strategy used for estimating the effectiveness of the selected supervision in improving the network’s domain generalization ability, two training datasets are used for cross-validation for the mobility-based part segmentation task, namely the training dataset and the auxiliary training dataset. While three training datasets are used for both primitive fitting and the semantic-based part segmentation. Three clusters out of four clusters are used for cross-validating in primitive fitting task. Three categories, including “Chair”, “Lamp”, and “StorageFurniture”, are used for cross-validating semantic-based part segmentation.

As for the whole cross-validating process, we train the network on each train-validation split fold for one

single epoch with intermediate supervisions added based on the selected supervision feature and the segmentation task related supervision optimized simultaneously. After that, the average metric gap across all train-validation splits is taken as the estimated generalization score for the selected supervision feature, which is also used as the reward score for further supervision feature distribution space update.

Ablation Study. For details of different ablated versions w.r.t. the supervision space design. “Less operants” denotes using a smaller input feature candidate set that is not expanded. For a set of input part-level matrix candidates we used in the full regular searching process where the full version of the input feature set that is expanded from input features, only the set containing part and geometry-aware matrices formed from input features directly is used in this version. More specifically, in “Less operants” setting, we use $\{P_i, N_i\}$ as the input feature set, while the full version contains 7 matrices. “Less unary operators” denotes the version where only a subset of unary operators is used as the unary operator set. Compared with the full version listed in Table 6, the one used in the ablated version is {Identity, Square, Double, Negative}. “Less binary operators” refers to the version where only a subset of binary operators is used as the binary operator set. Compared with the full version listed in Table 6, the one used in the ablated version is {Add, Minus, Multiply}. “Less tree height” means the maximum height of the operation tree, which is measured by the number of the connected operation cells. Compared with the one used in the full version that is set to 3, the maximum tree height in the ablated version is set to 2.

Models for comparison with HPNet. In the primitive fitting task, we develop two models to compare with HPNet fairly, considering the clustering-based network architecture used in HPNet that is different from the classification-based segmentation module used in our default setting to evaluate AutoGPart. Two models are designed by 1) replacing the clustering-based segmentation module used in HPNet with a classification-based segmentation module, denoted as “HPNet*”; 2) plug AutoGPart in the first learning stage of HPNet to search for useful intermediate supervisions that can help the network learn representations using more invariant features and avoid using shortcut features [8], denoted as “AutoGPart_{HPNet}”. Following are some detailed settings for such two models. For HPNet*, the network is formed by replacing the clustering-based segmentation module with a classification-based segmentation module with supervisions added on the first learning stage kept. The model is trained and evaluated using the same setting as for our own model. That is, train the model for 200 epochs and select the best validation epoch for further evaluation. For AutoGPart_{HPNet}, we also adopt a two-stage training

procedure. In the first stage, AutoGPart is applied to search for useful intermediate supervisions using the gap of the contrastive-style loss between the training dataset and the validation dataset across all train-validation splits as the reward. After the supervision distributions have been optimized, we greedily select what supervisions to use and plug them in the first learning stage of HPNet by optimizing such losses and other losses introduced in HPNet simultaneously. The resulting model optimized in the first learning stage is then taken for further evaluation using the clustering-based segmentation module.

Point-cloud processing backbones. PointNet++ used in our default setting is the same one as that used in SPFN [24]. DGCNN used in the default setting is the same one as that used in HPNet [47] for representation learning.

Input features. For semantic-based part segmentation task where per-point normal vector is not contained in input features, we estimate a normal vector for each point using open3d [57]. For primitive fitting, input geometric features for each point i contain a coordinate vector \vec{p}_i and a ground-truth normal vector \vec{n}_i . For mobility-based part segmentation task, input features for each point i are composed of a coordinate vector \vec{p}_i and a flow vector \vec{f}_i . The flow vector is estimated according to two input shapes.

Baselines. For domain-agnostic baselines for general domain generalization problems, like MixStyle [56], Meta-learning [21], Gradient Surgery [31], we implement them for segmentation tasks carefully with reference to their released code. For mobility-based part segmentation task, Deep Part Induction [48] is a learning based segmentation network. Although the test dataset used in our model is the same as the one used in their work, we download the code, re-implement it using PyTorch, and further train it using our training dataset where each pair is generated on-the-fly from inferred meta-data for part mobility information. Note that the performance reported in the original paper (77.3% MIOU on the test set) is achieved using several iterations between flow estimation and part segmentation. The performance of the model using a single iteration is 63.1% as reported in their work. It is also different from the one we report, probably due to the different training dataset. The performance of other baselines such as JLinkage clustering (JLC) [51] and (Spectral Clustering) SC [42] are taken from [48] due to the same test dataset.

For primitive fitting, SPFN [24] and HPNet [47] are two task-specific methods to solve the problem. We download the code of SPFN released by the author, carefully re-implement it using PyTorch and test the model on the same train-validation-test split as the one used for our model. For HPNet, we download the official implementation and test the model’s performance on our train-validation-test split.

For the semantic-based part segmentation, Learning to Group [30] is a two-stage learning-based segmentation network with a representation learning stage and a reinforcement learning based strategy for part segmentation; SGPN [43] and GSPN [50] are also two task-specific segmentation strategies. WCSeg [17] is a traditional segmentation method. The performance of those methods are directly taken from [30] due to the same training and test dataset.

Software configurations. We use Python 3.8.8 and PyTorch 1.9.1 to write the main code framework. Other main packages used include torch_cluster 1.5.9, torch_scatter 2.0.7, horovod 0.23.0 for PyTorch, etc.

Hardware configuration. All training experiments, including supervision search stage and regular training stage, are conducted on 8 NVIDIA Geforce RTX 3090 GPUs in parallel. Experiments for inference stage is conducted on one single NVIDIA Geforce RTX 3090 GPU.

D. Additional Experimental Results

D.1. Contrastive Learning based Part Segmentation

Our main experimental results have proved the effectiveness of HPNet [47] on the primitive fitting task. It is a carefully designed two-stage framework with a representation learning network that learns feature representations, parameters and other geometric features such as normal vectors and further hybrid such learned features for the following Mean-Shift clustering module. Such design achieves impressive performance on the primitive fitting task, with MIoU performance 79.5%. It is not clear whether such network architecture is suitable for other two tasks. To apply them on the mobility-based part segmentation and semantic-based part segmentation, we abbreviate its first representation learning stage to only learn per-point representations optimized by a contrastive style loss and further fit the optimized representation to the clustering stage.

We conduct experiments by applying such contrastive learning based part segmentation strategy to the mobility-based part segmentation task and the semantic-based part segmentation task using PointNet++ [36] as the backbone. Results are summarized in Table 8. It can be inferred that although such contrastive based segmentation networks can work well on the primitive fitting task, it cannot get satisfactory results on the mobility-based part segmentation task or semantic-based part segmentation task using PointNet++ as the backbone. Possible reasons may include 1) Such network architecture is not a universal one for all segmentation tasks, considering that it is originally designed for the primitive fitting task. 2) This network architecture is not universally suitable

for both PointNet++ and DGCNN. For comparison, our strategy is more universal compared with the contrastive learning based segmentation strategy.

Table 8. The performance of the contrastive learning based part segmentation networks (“Contrastive”) on the mobility-based part segmentation task (“Mobility”) and the semantic-based part segmentation task (“Semantic”). “PN++” refers to “PointNet++”. Reported values are performance of the trained networks on out-of-domain test datasets.

	Mobility	Semantic
Contrastive	44.0	25.9
AutoGPart _{PN++}	66.5	33.9

D.2. Part-aware and Geometry-discriminative Supervision Feature Space

In our method, we design the supervision space to be aware of part-level information by sampling features from points having the same part labels with the target point. Moreover, the feature space is made aware of geometric features by using ground-truth geometric features to construct the input feature set. In this section, we conduct ablation study on the mobility-based part segmentation task to prove the necessity of making the calculated features aware of both part labels and geometric features.

Ablating part labels. We conduct an experiment to ablate part labels in the calculation of supervision features by using features of points sampled from the neighbourhood of each point without considering whether they belong to the same part. Resulting models are denoted as “–Part Labels” in Table 9.

Ablating geometric features. We conduct an experiment to ablate geometric features in the supervision feature calculation process by only using part labels from the sampled points as input features in the supervision feature space. Resulting models are denoted as “–Geometric Features” in Table 9.

The results of such ablations are summarized in Table 9. Such results demonstrate the necessity of including geometric features and part labels in the intermediate supervision feature calculation process simultaneously. Furthermore, being aware of geometric features are more important than including part labels in the supervision feature calculation process if we compare the performance of the resulting models for their respective ablation versions. Moreover, if we compare the performance of the model “–Geometric Features” with the model optimized using no intermediate supervisions, it can be discovered that the ablated version, where only part labels are involved in the calculation process, would result very limited performance improvement. Possible reasons may include 1) Operators designed for geometric features are not that suitable to encode part labels; 2) Adding

Table 9. Ablation study w.r.t. reward function design and supervision space design. For abbreviations used, “Arch.” refers to “Architecture”; “PN++” denotes “PointNet++”; In/Out-of-dist. refers to In/Out-of-distribution performance.

Ablation	Arch.	In-dist.	Out-of-dist.
/	DGCNN	83.1	73.8
–Part Labels		82.4	71.3
–Geometric Features		83.9	68.9

Table 10. Ablation study w.r.t. reward function design and supervision space design. For abbreviations used, “Arch.” refers to “Architecture”; “PN++” denotes “PointNet++”; “No Loss” indicates simple segmentation networks trained without intermediate supervisions; In/Out-of-dist. refers to In/Out-of-distribution performance.

Ablation	Arch.	In-dist.	Out-of-dist.
/	PN++	87.2	66.5
/(No Loss)		86.4	61.0
NN		86.5	62.5
/	DGCNN	83.1	73.8
/(No Loss)		88.3	68.1
NN		89.5	67.0

intermediate supervisions by only letting the model aware of part labels encoded in another form is not enough to help the model learn more part related cues. Letting the calculated supervision features aware of part-level geometric information can help the network learn useful cues defining parts for the part segmentation task, thus benefiting the networks’ performance better than only using a part of them.

D.3. Supervision Feature Generation Strategy

In this work, we propose to generating supervision features from an operation tree operating on input part-aware geometric features. No doubt that there are many other methods that can be used to generate such features from ground-truth input features for supervision. We explore one possibility by letting a network consume such input part-aware geometric features for intermediate ground-truth supervision features. The network is optimized together with the main network by gradient descent. Resulting models are denoted as “NN” and the results are summarized in Table 10. As can be inferred from the table, using a network to generate ground-truth features for prediction is not a wise choice, with the performance improvement that can be added on the original simple segmentation networks very limited. It indicate that, using such simple gradient descent strategy, a network cannot “learn” the correct way to generate high-quality ground-truth features for intermediate supervisions. This can also prove an effective strategy to use a computing graph-like operation tree for ground-truth features.

D.4. Experimental Results Details

In this section, we present some details of experimental results that are not presented in the main paper.

Performance comparison between HPNet and AutoGPart_{HPNet}. Performance comparison between HPNet and AutoGPart_{HPNet} on different data clusters is shown in Figure 7.

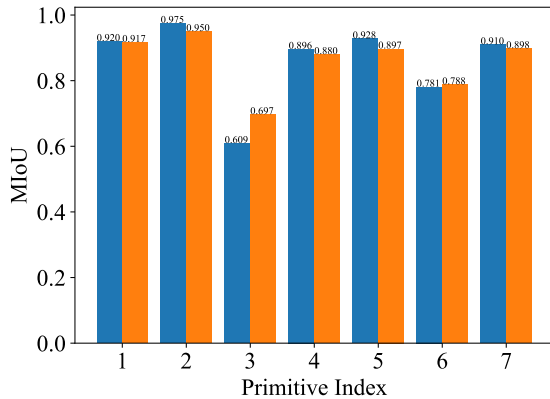


Figure 7. Performance comparison between HPNet and AutoGPart_{HPNet} on different data clusters.

D.5. Segmentation Results Visualization

In this section, we visualize and present some segmentation results for the mobility-based part segmentation task and the semantic-based part segmentation task. Point clouds are normalized into a ball with radius 1.0. Thus the presented shapes may be different from those used in the training stage.

Mobility-based part segmentation. Figure 8 and 9 (presented in next few pages) show the selected segmentation visualization for the mobility-based part segmentation task, where “Baseline” denotes “DGCNN” model using no intermediate supervision while “Ours” denotes “AutoGPart_{DGCNN}” model using intermediate loss searched by AutoGPart.

Semantic-based part segmentation. Figure 10 and 11 (presented in next few pages) show the segmentation results of AutoGPart on shapes from several test categories for the semantic-based part segmentation task.

E. Further Discussion

In this work, we propose to automatically find useful intermediate supervisions to help with improve the generalization ability of 3D part segmentation networks. Although experiments prove the value of adding intermediate supervisions for improving networks’

cross-domain performance, the network structure used in current work is still limited to a single stage end-to-end learning-based network which does not vary across different segmentation networks. However, it is valuable to explore how to include the architecture of segmentation networks into the design space and automatically select suitable network architectures for different part segmentation tasks. Making the network architecture flexible can help enlarge the design space, thus including more highly expressive networks into the search space.

Another line lies in more explorations on other tasks, not only limited to 3D part segmentation tasks. It is expected that adding intermediate supervisions a general approach to prevent network from learning shortcut features for tasks from a much broader range. However, it is still not proved and worth further exploration. Moreover, migrating the research goal from only improving the network's cross-domain generalization ability to enhancing both of its in-distribution and out-of-distribution performance is also a meaningful direction.

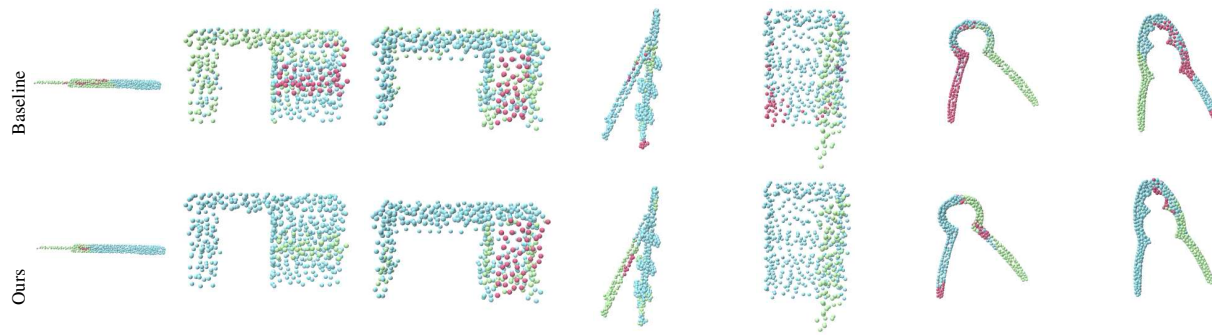
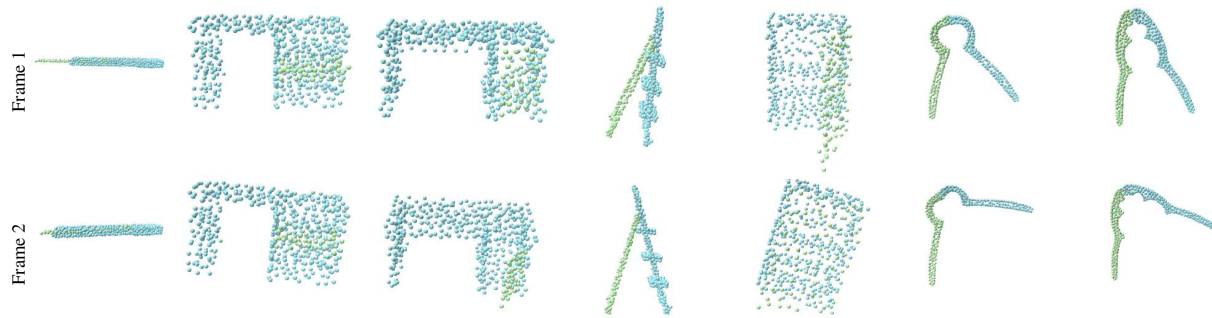
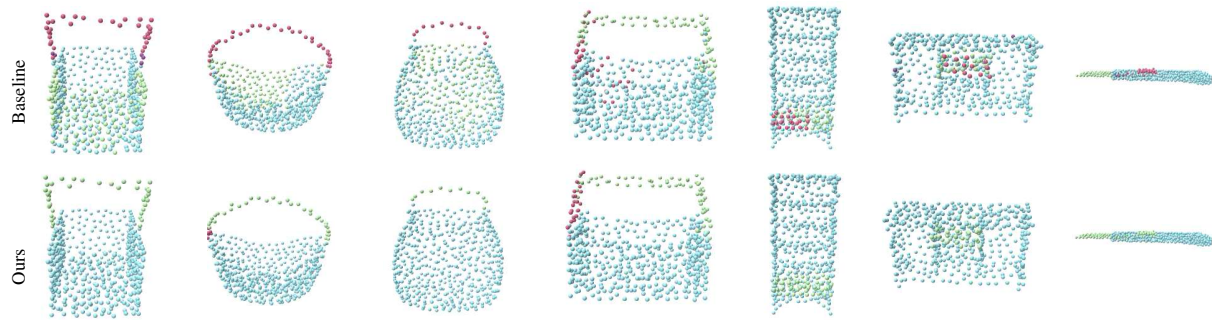
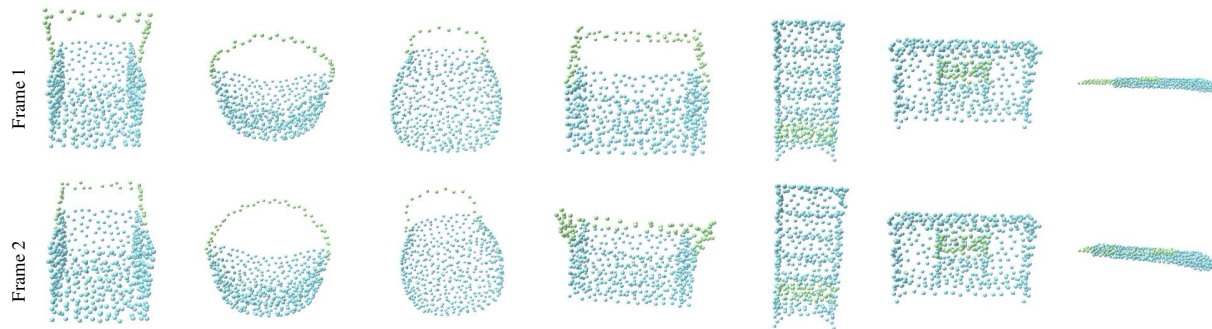


Figure 8. Mobility-based segmentation results visualization. “Baseline” refers to the segmentation network using DGCNN as its backbone without adding any intermediate supervision. “Ours” denotes the model “AutoGPart_{DGCNN}”.

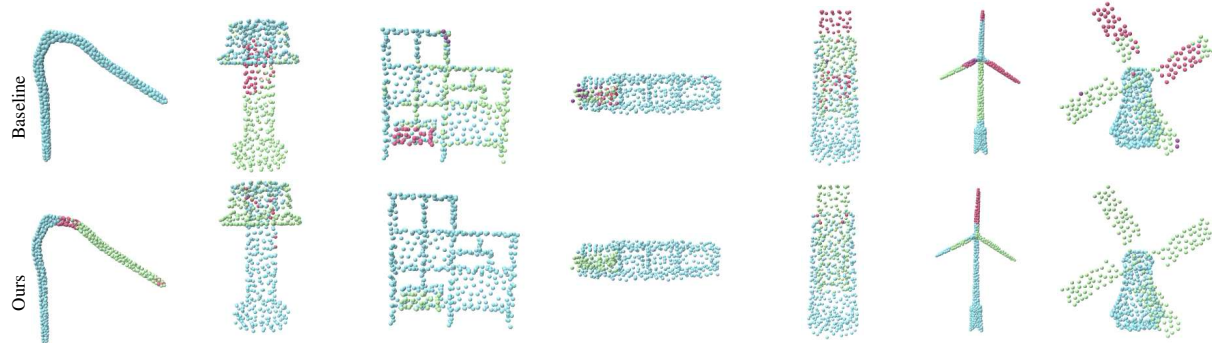
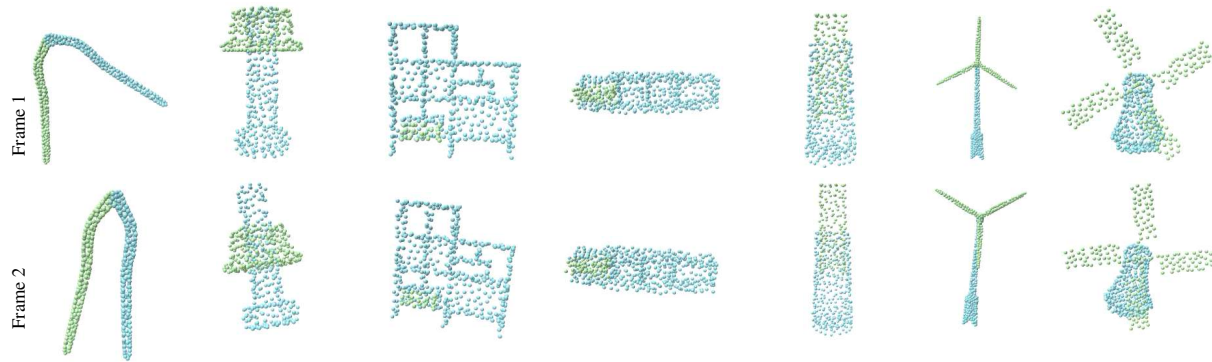


Figure 9. Mobility-based segmentation results visualization. “Baseline” refers to the segmentation network using DGCNN as its backbone without adding any intermediate supervision. “Ours” denotes the model “AutoGPart_{DGCNN}”.

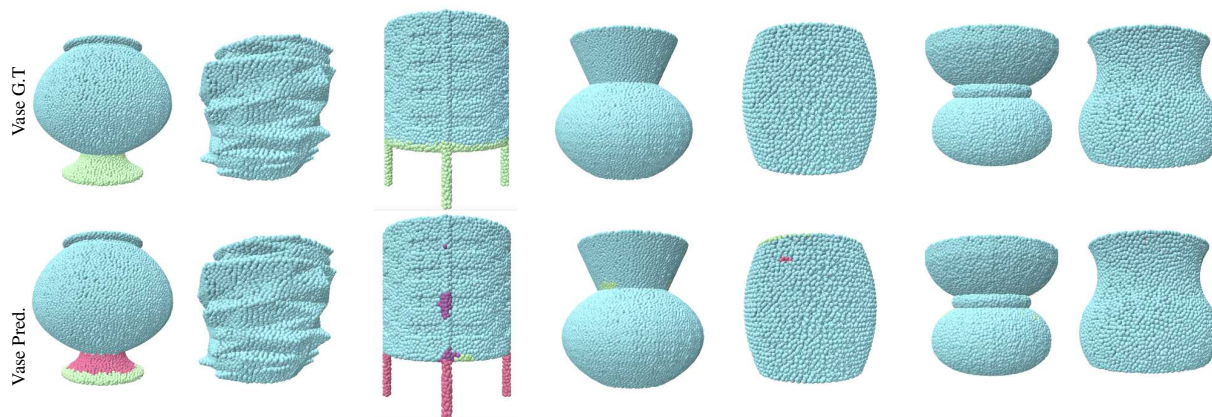


Figure 10. Semantic-based segmentation results visualization. “G.T” refers to the ground-truth segmentation results. “Pred.” denotes the model “AutoGPart”.



Figure 11. Semantic-based segmentation results visualization. “G.T” refers to the ground-truth segmentation results. “Pred.” denotes the model “AutoGPart”.