
ADVANCEMENTS IN POINT CLOUD-BASED 3D DEFECT DETECTION AND CLASSIFICATION FOR INDUSTRIAL SYSTEMS: A COMPREHENSIVE SURVEY

Anju Rani

Department of Energy, Aalborg University
Niels Bohrs Vej 8, Esbjerg, 6700, Denmark
aran@et.aau.dk

Daniel Ortiz-Arroyo

Department of Energy, Aalborg University
Niels Bohrs Vej 8, Esbjerg, 6700, Denmark
doa@et.aau.dk

Petar Durdevic

Department of Energy, Aalborg University
Niels Bohrs Vej 8, Esbjerg, 6700, Denmark
pd1@et.aau.dk

February 21, 2024

ABSTRACT

In recent years, 3D point clouds (PCs) have gained significant attention due to their diverse applications across various fields such as computer vision (CV), condition monitoring, virtual reality, robotics, autonomous driving etc. Deep learning (DL) has proven effective in leveraging 3D PCs to address various challenges previously encountered in 2D vision. However, the application of deep neural networks (DNN) to process 3D PCs presents its own set of challenges. To address these challenges, numerous methods have been proposed. This paper provides an in-depth review of recent advancements in DL-based condition monitoring (CM) using 3D PCs, with a specific focus on defect shape classification and segmentation within industrial applications for operational and maintenance purposes. Recognizing the crucial role of these aspects in industrial maintenance, the paper provides insightful observations that offer perspectives on the strengths and limitations of the reviewed DL-based PC processing methods. This synthesis of knowledge aims to contribute to the understanding and enhancement of CM processes, particularly within the framework of remaining useful life (RUL), in industrial systems.

Keywords Deep learning · Condition monitoring · Defect detection · Point cloud · Classification · Segmentation

1 Introduction

Condition monitoring (CM) is of vital importance in ensuring the longevity and proper maintenance of structures, such as bridges, buildings, industrial facilities, and infrastructure. Traditionally, visual inspection has been used over the years for CM applications. Traditional two-dimensional images face limitations in providing depth information and relative object positions, which is crucial for tasks involving spatial details such as autonomous driving, virtual reality, and robotics. The emergence of 3D acquisition technologies, including depth sensors and 3D scanners, has effectively addressed this limitation by facilitating the extraction of detailed 3D information. The utilization of 3D data offers a significantly improved understanding of objects compared to traditional 2D images. In recent years, there has been a growing emphasis among researchers on harnessing 3D scanned objects for defect detection and segmentation in industrial applications [1, 2, 3, 4]. The representation of 3D data can be of various forms, including depth images, PCs, meshes and volumetric grids. PC representation stands out for preserving the original geometric features in 3D space without any discretization, making it the preferred choice in many applications. The PC consists of unstructured 3D vectors, where each point represents a vector indicating its 3D coordinates (XYZ) with additional feature channels

such as colour (RGB values), intensity, and surface normals. Also, the PC exhibits properties like unstructured points, interaction among points, and invariance under transformation. These characteristics contribute to the flexibility and adaptability of PC representation in capturing complex geometric structures.

In the last decade, DL has emerged as the most influential technique in the field of 2D-CV such as image recognition, object detection, and segmentation. However, the application of DL to 3D PC data presents unique challenges due to the unstructured, high-dimensional, and disordered nature of PCs. Traditional convolutional networks designed for regular grids may not be directly applicable to PCs. Therefore, raw PC data is pre-processed to make it compatible with DL algorithms. This involves steps such as noise removal, data cleaning, down-sampling, and normalization to enhance and ensure data consistency. Later, various network architectures, including convolutional neural networks (CNNs) [5, 6, 7], graph neural networks (GNNs) [8, 9, 10], or hybrid networks [11, 12, 13], can be used for specific tasks such as 3D classification and segmentation. The DL model is then trained using annotated PC data. This involves feeding the PC into the network, computing the loss between ground truth labels and predicted labels, and then updating model parameters through back-propagation. The training stage often requires a large input dataset, prompting the use of data augmentation techniques to improve generalization. After training and evaluation, the model can be used for inference on new, unseen PC data, followed by post-processing steps to refine the model output. A taxonomy of existing DL methods for processing 3D PCs is shown in Figure 1.

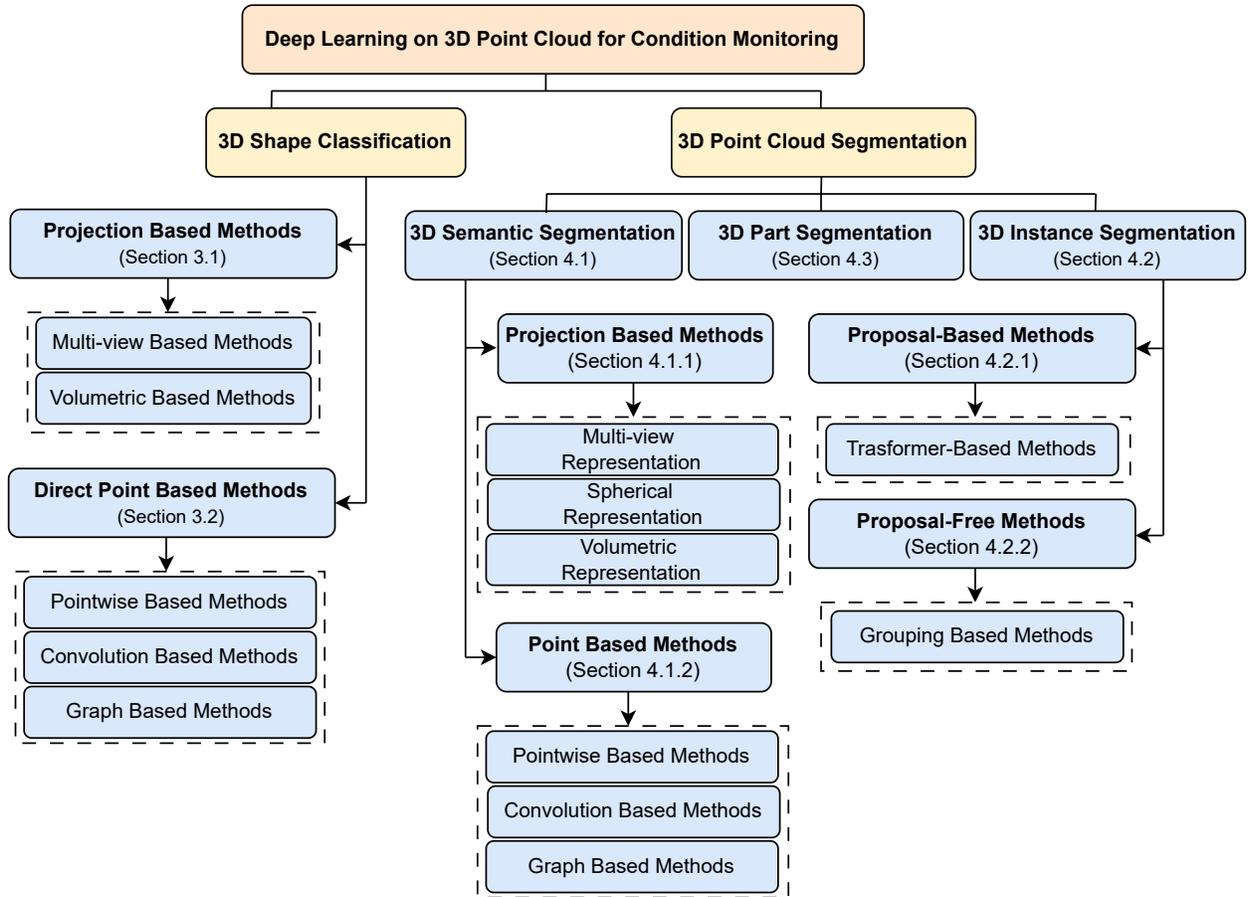


Figure 1: A taxonomy of DL methods for processing 3D PC data.

The paper provides a comprehensive review of DL methods applied to 3D PC data, with a specific emphasis on their applications in industrial settings. While previous reviews have explored DL techniques using standard datasets, this paper goes beyond by dissecting fundamental methodologies and recent advancements in 3D shape classification and segmentation, specifically catering to CM requirements in industrial applications. The review covers both traditional and innovative approaches, shedding light on the inherent challenges and potential solutions in processing 3D PC data for CM applications in industrial settings. Additionally, it provides a detailed summary of existing DL methodologies for feature learning in 3D PCs, outlining their respective strengths and weaknesses. The inclusion of publicly available datasets relevant to 3D shape classification and object segmentation enhances the practical value of the discussion.

Overall, the synthesis of existing knowledge in this review aims to identify gaps in the current understanding and pave the way for further innovations in the dynamic field of 3D PC data processing, offering valuable insights for researchers and practitioners. The key contributions of this review paper encompass the following aspects:

1. The paper provides a thorough survey of the most recent advancements in DL-based 3D PCs applied to both traditional and CM applications. The discussion is categorized into two main domains—shape classification and 3D object segmentation.
2. The review systematically compares and summarizes recent methods for CM, with a specific focus on damage detection in industrial applications. This comparative analysis not only highlights the diverse approaches but also provides an insightful assessment of the strengths and limitations of each method, offering valuable guidance for researchers and practitioners.
3. The paper goes beyond the current state of the field by offering valuable insights into potential future research directions and applications in the realm of deep learning-based condition monitoring using 3D PCs. This forward-looking perspective aims to inspire and guide future research endeavours in the dynamic and evolving field.

The structure of this review paper is organized as follows: Section 2 discusses the existing datasets and evaluation metrics utilized for 3D PC classification and segmentation tasks. Section 3 focuses on DL methods used for 3D shape classification, unravelling the evolution and applications of these methodologies. In Section 4, an extensive survey is conducted on existing methods for 3D PC segmentation, including semantic segmentation, instance segmentation, and part segmentation. The review concludes in Section 5, synthesizing insights and outlining future research directions.

Table 1: Available benchmark PC dataset for classification and segmentation.

Ref.	Dataset	Description	Year	Classes	Object/Point Count	Classification	Segmentation
[14]	Oakland	Urban environment	2009	44	1.6 M	✓	
[15]	ISPRS	Buildings, trees, and 3D building reconstruction	2012	9	1.2 M		✓
[16]	Paris-rue-Madame	Street in Paris	2014	17	20 M	✓	✓
[17]	IQmulus	Dense urban environments	2015	22	300 M	✓	✓
[18]	ScanNet	Indoor Scenes	2017	20	2.5 M	✓	✓
[19]	S3DIS	Structural elements	2017	13	273 M		✓
[20]	Semantic3D	Robotics, augmented reality and urban planning	2017	9	4000 M	✓	
[21]	Paris-Lille-3D	Objects in urban environment	2018	50	143 M	✓	
[22]	SematicKITTI	Autonomous driving	2019	28	4549 M		✓
[23]	Toronto	Urban roadways	2020	9	78.3 M		✓
[24]	DALES	Aerial geographical scan	2020	9	505 M	✓	✓
[25]	nuScenes	Autonomous driving	2020	7	5 B	✓	✓
[26]	ModelNet	CAD-generated objects	2015	662	1.3 M	✓	
[27]	ShapeNet	CAD-generated objects	2015	3,135	300 M	✓	
[28]	ModelNet40-C	Corruption robustness	2022	40	1.85 M	✓	
[29]	ScanObjectNN	Scanned indoor scenes	2019	15	15,000	✓	
[30]	STPLS3D	Synthetic and real aerial photogrammetry	2022	20	15,888		✓
[31]	SUN RGB-D	3D room layout and scenes	2015	700	10,335		✓
[32]	Hypersim	Synthetic indoor images	2021	461	77,400		✓

2 Background

2.1 3D Datasets

The availability of publicly accessible datasets plays a pivotal role in facilitating the analysis and comparison of various models in the domain of 3D PC applications. Researchers have curated diverse datasets specifically designed for tasks such as 3D shape classification, 3D object detection, and 3D PC segmentation. Table 1 provides a concise summary of these benchmark datasets along with their descriptions. These datasets can be broadly categorized into two main types: real-world and synthetic datasets. In real-world datasets [18, 29], the objects are occluded at varying levels while some objects may contain background noise. On the other hand, objects in synthetic datasets [26, 27], are without any occlusion and background noise, offering a controlled environment for experimentation.

2.2 Evaluation Metrics

Different evaluation metrics are employed in the literature to assess the performance of deep learning-based 3D PC processing tasks. For 3D shape classification, the most common performance criteria include *overall accuracy (OA)* and *mean class accuracy (mAcc)* respectively. *OA* represents the mean accuracy for all test instances while *mAcc* is the mean accuracy for all shape classes. In the case of 3D PC segmentation, *OA*, *mAcc*, *mean intersection over union (mIoU)* and *mean average precision (mAP)* are the most frequently used performance criteria. *OA* in this case represents the mean accuracy for PC segmentation, *mAcc* depicts the mean accuracy for different classes in segmentation and *mIoU* Measures the overlap between predicted and ground truth segments. Particularly, *mAP* is used in instance segmentation of 3D PCs. These metrics provide a quantitative assessment of the performance of deep learning models across various 3D PC processing tasks. However, the appropriate metric is chosen based on the specific task and the desired aspects of performance to be evaluated.

3 Deep Learning for 3D Shape Classification

The existing 3D shape classification methods can be broadly categorized into two major groups: projection-based methods and direct point-based methods. Figure 2 depicts various milestone methods within these categories, showcasing the diversity of approaches discussed in the literature.

3.1 Projection-Based Methods

These methods typically involve the projection of a 3D PC into 2D images, facilitating the application of well-established 2D image processing techniques for classification tasks. This category encompasses techniques that leverage multi-view images or volumetric images to represent and analyze 3D shapes.

3.1.1 Multi-View Based Methods

This method captures 3D shape projections from multiple viewpoints and extracts features independently from each view. Traditional methods, such as CNNs, can be applied to each view to extract distinctive features, which are subsequently fused to classify the shape accurately. However, the effectiveness of methods largely depends on the number of views selected for the classification. Multi-view CNN (MVCNN) [33] captures 3D shapes from various viewpoints and passes them through CNNs to extract features independently from each view. These extracted features undergo max-pooling and are then passed through another CNN layer to generate a compact shape descriptor. However, max-pooling retains only the largest elements from the viewpoints, resulting in a loss of information. [34] implemented MVCNN for classifying ten-defects in road infrastructure. The author compared the classification performance between MVCNN and PointNet [35]. The results demonstrated the superior performance of MVCNN with mAcc of 0.98 in comparison to 0.83 in the case of PointNet. [36] proposed a CNN model to extract global features from regularly structured depth images. This approach contrasts with existing methods like MVCNN and PointNet, which utilize unstructured point cloud data. The depth images utilized in this study do not introduce any geometry loss, enabling fine-grid shape classification of defects in solder joints. Group-View CNN (GVCNN) [5] introduces a hierarchical shape descriptor by incorporating grouping and individual viewpoints information in the pooling process. While GVCNN exhibits a significant improvement in accuracy compared to MVCNN, it faces challenges, particularly with smaller views. Multi-view harmonized bi-linear network (MHBN) [37] combines local convolutional features from multiple views using bi-linear pooling to generate a global shape descriptor. Later, the sequential behaviour of the captured views was explored to recognise the 3D shapes. [38] combined CNNs and long short-term memory (LSTM) to aggregate multi-view features into shape descriptors. This approach leverages the temporal dependencies among views, enhancing the understanding of 3D shapes through the fusion of both spatial and sequential information. SeqViews2SeqLabels [39] takes into account the spatial relationship among viewpoints by introducing an encoder to aggregate the information from sequential views and a decoder for predicting global features or sequence labels. Subsequently, the author extends this approach with 3D2SeqViews [40], which efficiently aggregates information from both views and sequential spatial views in a hierarchical attention (view-level and class-level) mechanism. However, these methods are limited to aggregating ordered views and do not handle the aggregation of unordered views. Another hierarchical network based on view graph representation was introduced in view-based graph convolutional network (view-GCN) [41]. In this approach, the author constructed a view graph where multiple views are treated as nodes of the graph. The view-GCN learns discriminative shape descriptors based on the relationship between multiple views. Based upon this concept, multi-view GCN [42] was proposed for classifying defects (scratch, dent, protrusion) in synthetically generated 3D PC datasets on an aircraft fuselage. The author demonstrates the applicability of graph-based representations for capturing complex relationships among multiple views in the context of defect classification. Multi-view-based fusion pooling (MHFP) [43] adopts a hierarchical approach to fuse multi-view features into a compact descriptor, leveraging

correlations between several views. This method effectively removes redundant information while retaining maximum relevant information by using a 3D attention module to construct a graph. On a similar note, multi-view softpool attention networks (MVMSAN) [44] refines view feature information using a soft-pool attention convolution framework. The attention mechanism plays a crucial role in addressing challenges related to down-sampling, feature information loss, and insufficient detail feature extraction, ultimately contributing to improved performance of the model. With the recent success in vision transformer (ViT) [11], [12] proposed multi-view convolutional ViT (MVCVT), which combines CNN on each view to extract multi-scale local information and utilizes transformers to capture the relevance of multi-scale information across different views. This integration showcases the adaptability and effectiveness of transformer-based architectures in the context of multi-view feature extraction for 3D shape classification.

In summary, view-based methods learn from view features to obtain global feature descriptors, leveraging established CNN frameworks. Nevertheless, these methods often rely on traditional pooled downsampling techniques that prioritize retaining essential information rather than preserving the entirety of the data. This approach can result in insufficient extraction of view refinement feature information, leading to a substantial loss of valuable insights from the view features.

3.1.2 Volumetric-Based Methods

This method represents 3D shapes in the form of a 3D voxel grid using 3D volumetric convolutions such that each voxel signifies whether a point in 3D space is occupied by an object or not. VoxNet [45] addresses large PC data by integrating a volumetric occupancy grid with 3D CNN. [26] proposed a convolutional deep belief network, 3D ShapeNets to represent a 3D shape based on the probability distribution (PD) of binary variables on voxel grids. However, these methods do not perform well in processing dense 3D data due to high computation and memory requirements for higher resolution (computational complexity is a cubic function of voxel grid resolution) [46]. To overcome this limitation, a hierarchical compact structure needs to be introduced. OctNet [47] achieves this by partitioning 3D PC data hierarchically using a set of unbalanced octrees, where each leaf node stores a pooled summary of the features of the voxels. This approach focuses memory allocations on the relevant regions, enabling the use of deeper networks with high resolution. Subsequently, Octree-based CNN (O-CNN) [6] was proposed for 3D shape classification. O-CNN averages the normal vectors of a 3D model into fine-leaf octants as network input and performs 3D CNN over the octants occupied by the 3D shape surface. Another network based on the non-uniform indexing named Kd-Net [48] was introduced to mimic the convolutional-based network. Kd-Net requires small memory and computation in comparison to uniform grids. [49] used 3D grids to represent PC data, further expressed using 3D modified Fisher vector method. This vector acts as an input to the 3D CNN to produce global features. PointGrid [50], a hybrid network that integrates both the grid and point representation for efficient processing of the PC data. In [51], a multi-orientation volumetric DNN (MV-DNN), was proposed to limit the octree partition to a certain depth for reserving leaf octants with sparse features. This method improves classification for both low and high-resolution grids.

In summary, volumetric-based methods represent 3D PCs using voxel grids to address the unordered structure of the data. However, this approach requires input voxels to be in a regular form for convolutional operations, leading to information loss with low-resolution voxels and subsequently lower classification accuracy. Additionally, these methods face challenges related to high computation requirements, especially for high-resolution data.

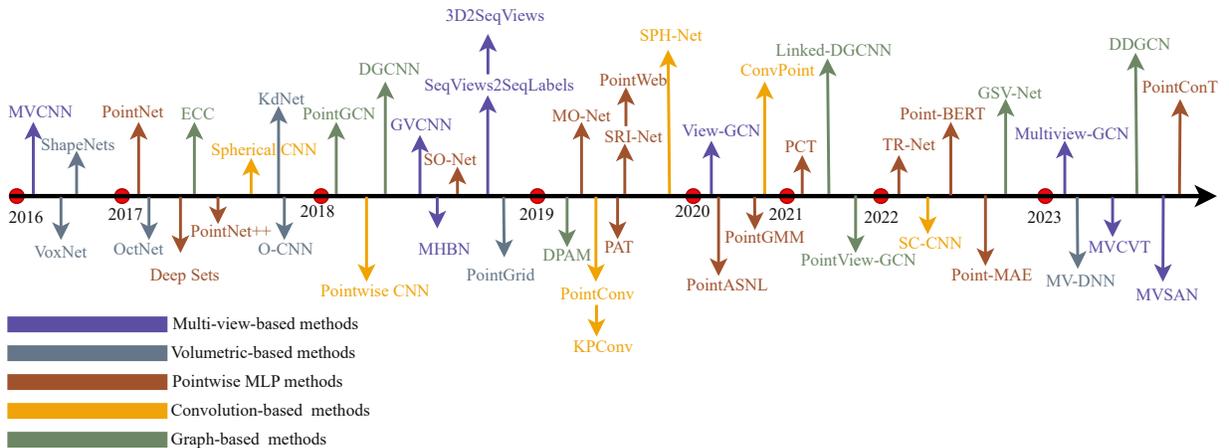


Figure 2: Chronological overview of the most relevant DL-based 3D shape classification methods.

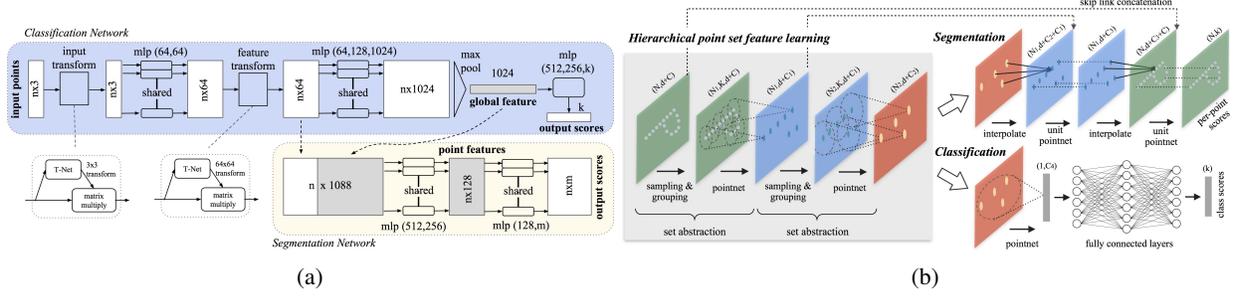


Figure 3: Architecture comparison of state-of-the-art methods; (a) PointNet [35] and (b) PointNet++ [52] respectively. Recreated from [35, 52]

3.2 Direct Point Based Methods

Direct point-based methods directly process the input PC data to produce a sparse representation. These methods extract a feature vector for each point by aggregating the features of neighbouring points. In this way, models designed for raw PC data typically begin by extracting low-dimensional features from individual points and later aggregate them to obtain high-dimensional features. Direct point-based methods can be further categorized into point-wise multi-layer perceptron (MLP), convolution-based, and graph-based methods.

3.2.1 Pointwise MLP Methods

These methods process each point independently through shared MLPs to extract local features. These local features are later aggregated to obtain global features using a symmetric aggregation function. PointNet [35] model represents unordered PCs as a set of 3D points, extracting local features independently for each point through multiple MLP layers. Global features are then obtained through max-pooling layers. Building upon this, the PointNet++ [52] model introduces a hierarchical structure incorporating layers such as the sampling layer, grouping layer, and PointNet learning layer (set abstraction level) to capture finer geometrical structures from neighbouring points. Figure 3 provides a visual comparison of the architectures of the PointNet [35] and PointNet++ [52] models respectively. [53] investigated the use of PointNet to detect defects (scaling, delaminations, and spalls) on the bridge surfaces. Due to the simple and efficient network, various models based on PointNet have been proposed in the literature for the direct processing of 3D PC data. [54] introduced a dual-level-defect detection PointNet (D3PointNet) for inspecting defects, specifically solder paste patterns in printers, through segmentation and multi-label classification. The author defined two hand-crafted features namely, edge and prior features to prevent loss in spatial information of the PC during processing. Self-organizing networks (SO-Net) [55] achieves permutation invariant for unordered PCs by building a self-organising map based on the spatial distribution of PCs. The hierarchical feature extraction of SO-Net results in a single feature vector that represents the entire PC. To enhance the performance of PointNet++, [56] proposed PointNetXt, introducing an inverted residual bottleneck design with separable MLPs into the PointNet++ architecture. This modification results in an effective and efficient model with a $10\times$ faster inference.

Several networks in the literature leverage geometrical features for 3D point cloud processing. Based on PointNet [35], Motion-based network (MO-Net) [57] incorporates the context of 3D geometry in the form of a finite set of moments as network input. This approach uses an attention mechanism to learn fine-grained local features of the PC. Point attention transformers (PATs) [58] represent each point in the PC using its absolute and relative positions concerning its neighbours. Then, group shuffle attention (GSA) captures the relations between these points, and a differentiable, permutation invariant, and trainable end-to-end gumbel subset sampling (GSS) layer is developed to learn hierarchical features. PointNet++ [52] based networks, such as PointWeb [59], explore the interaction between points using an adaptive feature adjustment module. This module interconnects all point pairs in a local region forming a fully-linked web to describe local regions for 3D recognition. However, these methods require large sample sizes that might not be available in real manufacturing applications. To address this limitation, [60] proposed a tensor voting-based approach to classify surface anomalies. Tensor voting makes inferences on the geometrical information such as curvature, surface and junction via voting over the neighbourhood for selecting the points which may include potential anomalies. Aggregated descriptive features are used for each selected PC sample and then passed to a sparse multi-class SVM classifier for anomaly classification and feature selection. Strictly rotation-invariant network (SRI-Net) [61] projects the PC data into a rotationally invariant representation, utilizing a PointNet backbone to extract global features and a graph aggregation method to extract local features. PointASNL [62] adaptively adjusts the coordinates of the initially sampled points using the furthest point sampling (FPS) algorithm and introduces a local-on local (L-NL) module to

capture local and long-range dependencies of these sampled points. PointGMM [63] is a coarse-to-fine feature learning method that subdivides the input point data into distinct groups using a hierarchical Gaussian mixture model (hGMM). This approach focuses on learning features of small and large regions, respectively. Here bottom GMM focuses on learning features of small regions while top GMM learn features of the larger regions respectively. [64] proposed a novel improved PointNet++ for classifying and segmenting the defects of different shapes and sizes in the sewer pipes. The author improved the network structure by incorporating residual connection and cross-entropy loss with label smoothing in the network. Later, the training process is optimized by AdamW and cosine learning rate decay.

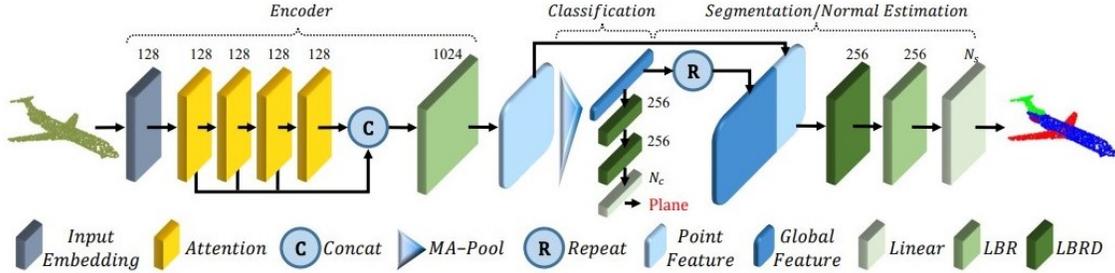


Figure 4: Architecture of PCT [65]. The encoder layer consists of an input embedding module with a stacked attention module while the decoder contains multiple Linear layers. Recreated from [65]

Unlike the above methods, PC transformer (PCT) [65] is based on a permutation-invariant transformer rather than a self-attention mechanism for handling unstructured and disordered point data with irregular domains. The overall architecture of PCT is presented in Figure 4. PCT transforms (encodes) the raw PC into a new feature space to characterize the semantic affinities between points. These features are then fed into the attention module to learn the discriminative representation for each point, followed by a linear layer to generate the final output feature [69]. Meanwhile, 3DMedPT [13] proposed a transformer-based network for analysing 3D medical PC data. Similarly, a Transformer-based network (TR-Net) [70] utilizes a neighbourhood embedding strategy and residual backbone with skip connections to enhance context-aware and spatial-aware features. The author uses an offset attention operator on PC spatial information to sharpen the attention weights for improving the extraction of global features. Inspired by the bi-directional encoder in transformers (BERT), Point-BERT [71] adopts a strategy of dividing the PC into distinct local blocks, generating discrete point labels that represent local information using a PC marker. This approach allows the model to capture specific details and features within localized regions of the PC. Similar to BERT, Point-BERT introduces a masking mechanism where some input PCs are randomly masked and then fed to the backbone transformer network. This facilitates bidirectional learning and enhances the model’s ability to capture contextual relationships in 3D PC data.

However, the uneven distribution of information in the PC may lead to a loss of information during the reconstruction task. To address this challenge, masked autoencoder (MAE) methods, such as Point-MAE [72] have been proposed. Point-MAE is a self-supervised learning (SSL) method designed to mitigate issues related to uneven information density and information leakage of PC locations. In a study by [68], a deep autoencoder network was proposed for processing 3D PC data of concrete bridges. The network takes encoded shape and neighbourhood features as inputs and uses a one-class support vector machine (OC-SVM) to classify spall defects on the concrete bridge’s PC data. The author tested the network on a diverse set of quasi-real PCs covering a variety of noise and defect conditions. PointConT [73] presents a novel approach to 3D PC processing by leveraging transformer-based clustering and self-attention mechanisms. The method focuses on clustering points based on their content and subsequently applying self-attention within each cluster. This design aims to capture long-range dependencies within the PC while managing computational efficiency. Additionally, the authors introduce an inception feature aggregation module, featuring a parallel structure to aggregate high and low-frequency information separately. [67] proposed a transformer-based PC classification network (TransPCNet) for detecting defects in sewer pipelines. TransPCNet comprises a feature embedding module responsible for extracting features from local neighbours, an attention module designed to learn and enhance feature extraction, and a classification module. Additionally, the authors introduced a weighted smoothing cross-entropy loss to aid the network in feature learning while addressing imbalances in the PCs.

In summary, pointwise MLP methods demonstrate efficiency and effectiveness in processing raw 3D PC data, leveraging simplicity to capture local features independently for each point, allowing for fine geometric structure understanding. Despite their advantages, challenges arise when handling large-scale and complex point clouds due to limitations in capturing long-range dependencies and holistic context. Additionally, these methods face difficulties in accommodating variations in point density, leading to potential impacts on the robustness of feature extraction. The introduction of point-based transformers and related models addresses some of these challenges by leveraging permutation-invariant

Table 2: Performance evaluation for classification methods on industrial applications.

Ref.	Application	Classes	Method	Results	Points/Objects
[66]	Defect classification in sewer	4 classes: normal, displacement, brick, and rubber ring	DGCNN	OA = 47.9	17,027
			PointNet	mIoU = 46.1 OA = 18.4 mIoU = 18.5	
[34]	Classification of infrastructure elements	10 classes: column, 3 types of culverts, 5 types of walls, and sump	PointNet	Mean F1 score = 89.3	1,496
			PointNet	OA = 83 F1 score = 87	
[60]	Classification of anomalies on steel surfaces	5 classes: debris, oscillation, slag, depressions, and pinholes marks	Tensor voting	Mean Acc = 86.27	96,266
[42]	Defect classification in precast concrete specimen	2 classes: defective, and normal	MVGCN	Euclidean = 97.9	2000
			DGCNN	Geodesic = 93.8 Euclidean = 70.8 Geodesic = 81.3	
[64]	Defect classification in concrete sewer pipes	5 classes: 3 circular defects of varying diameter, square and triangular defect	Improved PointNet++	Mean F1 score = 68.15	1.4 M
			PointNet++	Accuracy = 73.01 Mean F1 score = 61.36 Accuracy = 67.55	
[67]	Defect classification in polyvinyl chloride-sewer pipes	4 classes: normal, and defective (brick, rubber ring, displacement)	TransPCNet	F1 score = 60.58	17,027
			DGCNN	Precision = 61.47 F1 score = 16.66 Precision = 34.55	
			PointNet	F1 score = 30.23 Precision = 28.61	
[53]	Classification of surface defects on bridges	4 classes: cracks, spalling, scaling, and delaminations	PointNet	mAcc = 85.7	21 M
[68]	Spall Classification on bridges	2 classes: normal and defective	Point-wise	mAcc = 98	21 M
			PointNet	Precision = 68 mAcc = 97 Precision = 61	
[54]	Classification of printer defects	2 classes: normal and defective solder patterns	D3PointNet	mAcc = 97.17	4.2 M
			SDCNN	Precision = 97.28 mAcc = 98.1	800
[36]	Classification of solder joints shapes	2 classes: normal and defective	MVCNN	Precision = 83.9 mAcc = 93.6 Precision = 76.9	

transformers. These transformer-based approaches excel in managing unstructured and disordered point data, presenting a promising avenue for advancing the processing of 3D PC data.

3.2.2 Convolution-Based Methods

Following the remarkable success of CNNs in CV tasks [74, 75] such as image classification, object detection, and segmentation, there has been a significant effort to extend these methodologies to analyze geometric and spatial data. Unlike the regular grid structure in 2D images, geometric data (PCs, 3D models, etc.) lacks underlying grid information necessitating the development of new methods. Several convolution-based methods, including *continuous* and *discrete convolution-based methods*, have been developed for analyzing 3D PC data [76, 77, 78]. 3D continuous convolution methods are defined in a continuous space where weights for neighbouring points are spatially related to their centre point. Conversely, 3D discrete convolutions involve a fixed-size kernel sliding over a structured point grid, with weights assigned to neighbouring points determined by their offsets relative to the centre point of the kernel.

Among *continuous convolution methods*, PointConv [79] stands out by representing convolution kernels as nonlinear functions of the local coordinates of 3D points. These functions comprise weight (learned with MLP layers) and density (learned by kernel density estimation) functions. PointConv efficiently computes the weight function, providing translation and permutation-invariant convolution in 3D space. Another notable method, KPConv [80], introduces

a deformable convolution operator that learns local shifts at each convolution location, enabling adaptation of the kernel shape based on the input PC’s geometry. ConvPoint [81] takes a different approach by introducing a dense weighing function to define detailed and adaptive convolutional kernels. In this method, the derived kernel is explicitly represented by a set of points, each associated with specific weights. However, the above methods do not consider PC distribution while defining the convolution operator.

In the realm of *discrete convolution-based methods*, PointCNN [82] is a pioneering work that tackles the unordered and irregular structure of 3D PC data. PointCNN learns X-transformation from input PC data and permutes the weight of input point features into canonical order. CNN is then applied to these transformed features resolving the unordered and irregular structure of the 3D PC data. Pointwise CNN [83] applies the convolution operator on each point in the PC to learn pointwise features. The obtained outputs are then concatenated before being fed to the final convolution layers for segmentation or fully connected layers for object recognition. Unlike the traditional methods, Pointwise CNN does not require up-sampling or down-sampling of the PCs. Based on PointCNN, spherical harmonics network (SPH-Net) [84] proposed a rotation invariance CNN on PCs by using spherical harmonics-based kernels at different layers of the network. SC-CNN [85] implements a spatial coverage convolution by constructing an anisotropic spatial geometry in the local PC and replacing the depthwise convolution with the spatial coverage operator (SCOP). This method excels in learning high-order relations between points, providing shape information and enhancing network robustness.

In summary, continuous convolution methods, such as PointConv, KPConv, and ConvPoint, offer adaptability to diverse point cloud geometries and effective pattern capture. However, they overlook point cloud distribution considerations. In the discrete domain, PointCNN efficiently handles unordered structures, Pointwise CNN excels in pointwise feature learning, SPH-Net introduces rotation invariance, and SC-CNN learns high-order relations. While these methods enhance 3D PC analysis, challenges persist in distribution awareness and computational efficiency.

3.2.3 Graph Based Methods

Graph-based methods provide an alternative to CNNs for handling unstructured and unordered 3D point cloud data. Unlike CNNs, which operate on regular grid data, graph-based methods transform the point cloud into a comprehensive graph, avoiding the need for voxelization. A typical architecture of a graph-based PC network is illustrated in Figure 5. This approach allows for flexibility in capturing intricate relationships among points, representing each point in the point cloud as a vertex in the graph, with edges established between nearby points [86]. These edges analyze spatial relationships, creating a graph that encapsulates the geometric features of the original PC. Graph-CNN [87], also known as PointGCN classifies 3D PCs by combining localized graph convolution layers with two types of data-specific pooling layers (down-sampling). This method effectively incorporates the geometric information encoded in the graph, enhancing the robustness of the model. In contrast, Dynamic graph CNN (DGCNN) [8], inspired by PointNet addresses the limitation of processing each point independently, as in PointNet, leading to the neglect of local features between points. To solve this, Dynamic CNN uses the EdgeConv layer to capture edge features from each point and its neighbours. EdgeConv explicitly constructs a local graph while learning the embeddings for the edges, enabling the grouping of the points both in Euclidean and semantic space. [66] investigated the application of DGCNN and PointNet for classifying defects on synthetic and real sewer PC data. The author observed that the DGCNN network outperforms the PointNet network consistently for both synthetic and real datasets. Dynamic points agglomeration module (DPAM) [88] is based on graph convolution to agglomerate (sampling, grouping and pooling) points by multiplying the agglomeration matrix and points feature matrix. Based on PointNet and PointNet++, a hierarchical network is constructed by stacking multiple DPAMs by dynamically exploiting the relation between points and agglomerated points in a semantic space. Additionally, a variation of DGCNN, linked-DGCNN [89] simplifies the model by removing the transformation layer in DGCNN. This is implemented by connecting the hierarchical features of various dynamic graphs to address the issue of gradient vanishing. PointView-GCN [90] introduces a multi-level GCN to hierarchically aggregate shape features of single-view point clouds. This method allows the encoding of both object geometric cues and their multiview relationships, improving the extraction of global features. Gaussian super vector network (GSV-NET) [91] is a recent approach that captures and aggregates both local and global features of the 3D PC to enhance the information of the PC features. GSV-NET utilizes a combination of the GSV network and a 3D-wide inception CNN architecture to extract global features. The method then converts 3D point cloud regions into colour representations and employs a 2D-wide inception network to obtain local features. Also, [92] integrated the distance and direction in GCN (DDGCN) by constructing a dynamic neighbourhood graph. This dynamic graph utilizes MLPs and the similarity matrix to capture the local features of the PC. Additionally, the author modifies the loss function by incorporating the centre loss, enhancing the discriminative power of the model.

On the whole, point-based methods distinguish themselves by operating directly on raw PC data, rendering them well-suited for irregularly sampled and unstructured datasets with lower computational demands. Pointwise methods leverage MLP networks as fundamental building blocks for learning pointwise features, showcasing versatility in various network architectures. While literature indicates the superior performance of convolution-based networks for irregular

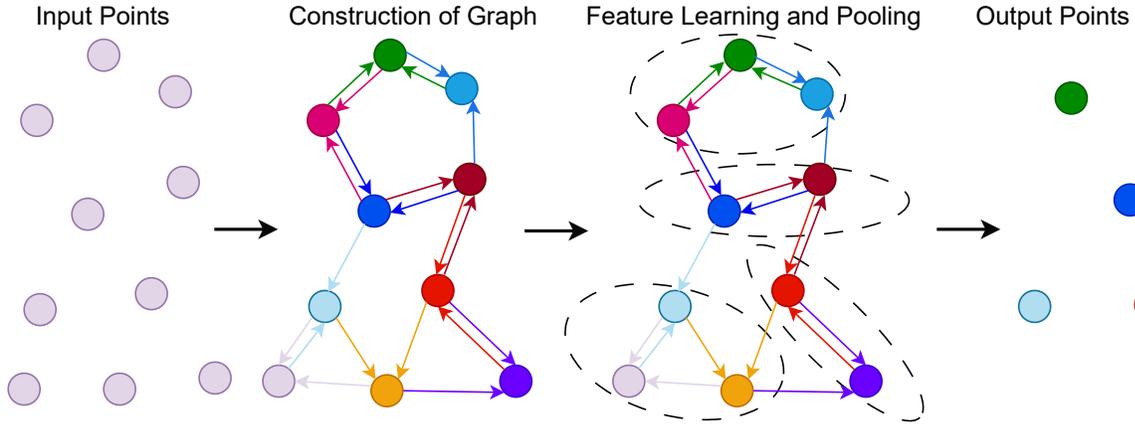


Figure 5: An illustration of a graph-based 3D PC network.

PC data, there exists limited research on both continuous and discrete convolution networks in this context. Graph-based approaches provide another avenue for handling irregular PC data, but extending these methods, particularly those based on spectral domain graph structures, to various graph configurations remains a challenging task. Future research directions may explore advancing convolutional and graph-based methodologies to enhance the understanding and processing capabilities of point-based methods for diverse and complex 3D datasets. Tables 2 present the outcomes of defect shape classification for industrial systems.

4 Deep Learning for 3D PC Segmentation

The task of 3D PC segmentation demands a comprehensive understanding of the geometric structure and intricate details of each point in the 3D PC data. The segmentation task can be broadly categorized into three major types:

1. Semantic segmentation (Scene level): This method classifies each point within a 3D PC into predefined categories by assigning semantic labels based on their characteristics, enabling a high-level understanding of the overall scene.
2. Instance segmentation (Object level): This method identifies and distinguishes each object in the 3D PC by assigning each point with a specific instance or object. Unlike semantic segmentation, which groups points into predefined categories, object level segmentation enables the recognition of separate instances of objects, even if they belong to the same semantic class.
3. Part segmentation (Part level): This method segments each component of the object in the 3D PC providing a more detailed object-level segmentation. Unlike semantic segmentation, which categorizes points into high-level classes, and instance segmentation, which identifies and distinguishes individual objects, part segmentation provides a more detailed breakdown of each object by segmenting its constituent parts.

These segmentation categories address different levels of abstraction, ranging from scene-level context to object-level identification and even detailed part-level segmentation. The annotated examples for semantic, instance and part segmentation on benchmark datasets are shown in Figure 7.

4.1 3D Semantic Segmentation

3D semantic segmentation, a key aspect of scene understanding, involves categorizing points in a 3D PC into predefined classes or labels. Similar to 3D shape classification, semantic segmentation methods can be divided into the following categories: projection-based methods (multi-view representation, spherical representation, and volumetric representation), direct point-based methods (pointwise MLP methods, convolution-based methods and graph-based methods) [2, 93, 94, 95]. Figure 6 illustrates the most recent methods in 3D semantic segmentation.

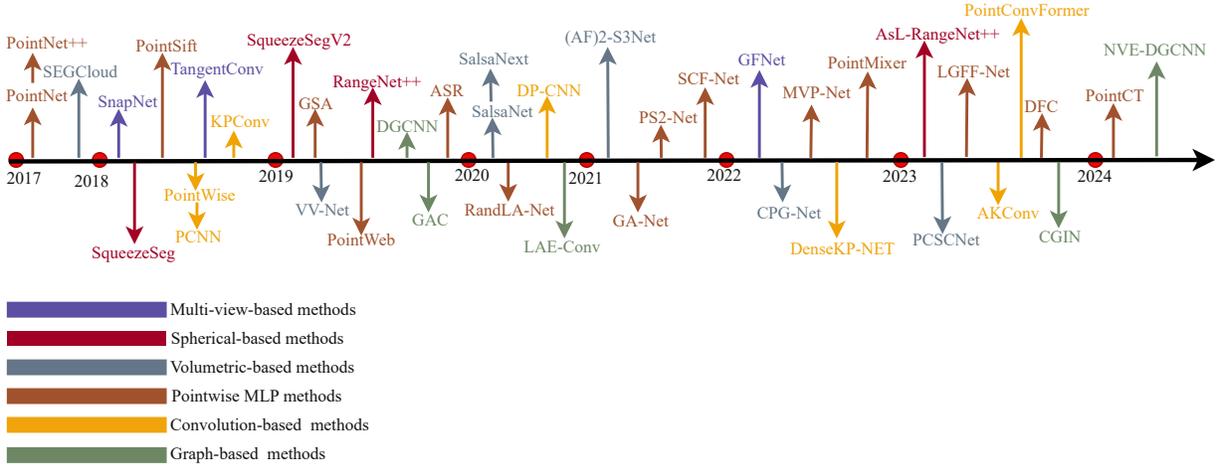


Figure 6: Chronological overview of the most relevant DL-based 3D semantic segmentation methods.

4.1.1 Projection-Based Methods

Projection-based methods in 3D semantic segmentation involve the transformation of 3D PC data into 2D images. This transformation is achieved through various techniques, including multiview projections, volumetric projections, and spherical projections.

Multi-View Representations: [96] projected PCs into 2D images using multiple camera views and then processed by fully convolutional networks (FCNs) for semantic segmentation. The resulting pixel-wise segmentation was re-projected into the original input PC. The final semantic label for each point is obtained by fusing the re-projected scores over the different views. However, there is a loss of information during the projection process. To address this limitation, [97] generated multiple RGB and depth (RGB-D) images containing geometric features using various camera positions. Pixel-wise labeling is then performed on these captured snapshots and fed to SegNet [98] and fusion is performed using residual correction [99] on the obtained predicted scores. SnapNet [100] selected specific snapshots of the PC to generate pairs of RGB-D images. Then pixel-wise labeling is performed on these 2D snapshots using FCNs. [101] proposed a novel tangent convolution for segmenting dense PCs. This method involves projecting local surface geometry onto a virtual tangent plane, serving as input for subsequent tangent convolutions. Generic flow network (GF-Net) [102] proposes a novel approach for learning geometric features by fusing information from multi-view representations. The author used KNN post-processing over KPConv to make it end-to-end trainable.

Overall, multi-view representation methods project 3D PCs into 2D images from various viewpoints for semantic segmentation. While providing diverse perspectives, they are sensitive to occlusions and viewpoint selection, impacting performance. Tangent Convolution addresses geometric information by projecting local surface features onto a virtual tangent plane. However, these methods may not fully exploit inherent 3D geometric information, leading to potential information loss.

Spherical Representations: Based on CNN, an end-to-end pipeline named SqueezeSeg [103] was proposed to provide labeled point-wise output data. This method utilizes a conditional random field (CRF) as a recurrent layer to further refine the segmented points. To reduce the impact of dropout noise on the accuracy of SqueezeSeg, the author proposed SqueezeSegV2 [104]. SqueezeSeg2 introduced a novel CNN module named context aggregation module (CAM) to aggregate contextual information from a large receptive field improving the network robustness to dropout noise. However, challenges persist in handling issues arising from intermediate representations, including blurry CNN outputs and discretization errors. RangeNet++ [105] overcomes these limitations by performing segmentation using CNN and an encoder-decoder hourglass-shaped architecture. The decoder incorporates a modified DarkNet [106] backbone architecture, enabling the use of aspect ratios beyond square configurations. Furthermore, RangeNet++ substitutes the CRF utilized in [103, 104] with GPU-based nearest neighbour calculations across the complete PC. However, when dealing with unbalanced training samples, the training outcomes may become skewed, leading to inaccuracies in segmentation results. The AsL-RangeNet++, an extension of RangeNet++, introduces an asymmetric loss (AsL) function proposed by [107]. This method uses the AsL function with Adam optimizer for calculating and adjusting object weights, enhancing the precision of semantic segmentation. However, these methods stack point data from

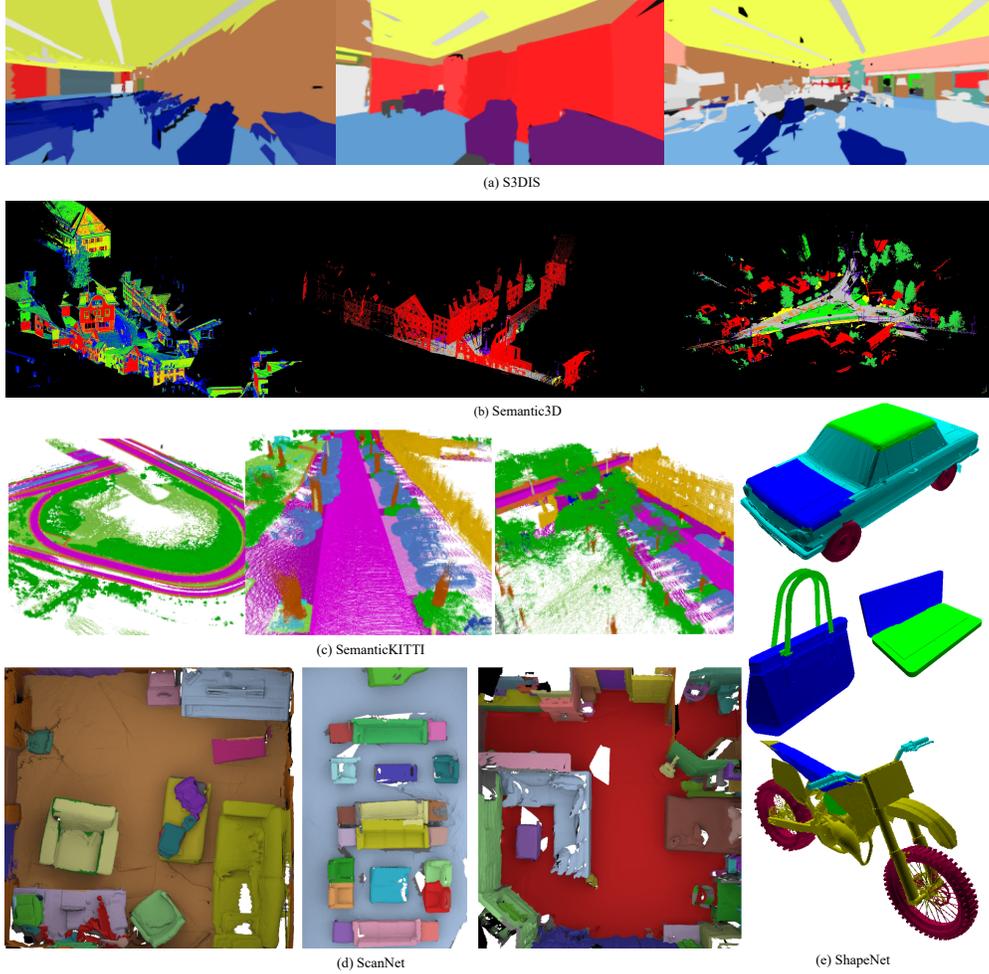


Figure 7: Annotated examples for (a) S3DIS [19], (b) Semantic3D [20], (c) SemanticKITTI [22] for 3D semantic segmentation, (d) ScanNet [18] for 3D-instance segmentation, and (e) ShapeNet [27] for 3D-part segmentation. Recreated from [18, 19, 20, 22, 27]

various modalities, such as coordinate, depth, and intensity, as inputs without accounting for their heterogeneous distributions.

Volumetric Representations: These methods transform unstructured 3D PCs into regular volumetric occupancy grids. The feature learning is then performed using NN to achieve semantic segmentation [27, 45, 108]. [109] projects the PC into occupancy voxels and fed into 3D-CNN to produce voxel-level labels, where all points within each voxel are assigned the same semantic label. [110] introduced InspectionNet, a 3D CNN-based framework designed to detect defects in synthetically generated concrete columns. SEGCloud [111] is an end-to-end framework for semantic segmentation that integrates NNs, tri-linear interpolation (TI), and fully connected CRF (FC-CRF). This approach generates coarse voxel predictions using 3D-CNN, which are then transferred back to the raw input 3D points through TI. Finally, FC-CRF is used to enforce global consistency and improve the semantic understanding of the points, resulting in fine-grained segmentation results. Voxel variational autoencoder network (VV-Net) [112] uses a combination of variational autoencoder (VAE) and 3D-CNNs to capture the point distribution within each voxel for semantic segmentation tasks. [113] introduced SalsaNet, an encoder-decoder network comprising a series of ResNet blocks in the encoder and employing upsampling and feature fusion in the decoder. Subsequently, SalsaNext [114] enhanced SalsaNet by replacing the ResNet encoder with a stack of residual dilated convolutions and a pixel-shuffle layer in the decoder, facilitating uncertainty-aware semantic segmentation. (AF)2-S3Net [115], an extension of S3Net [116] and S3CNet [117], is an encoder-decoder model designed for 3D semantic segmentation using sparse-CNN. In this approach, the encoder incorporates an attentive feature fusion module to capture both global and local features, while the decoder uses an adaptive feature selection module and feature map re-weighting to emphasize contextual

information obtained from the feature fusion module. Cascade point-grid fusion network (CPG-Net) [118], adopts a cascading approach to extract and aggregate semantic features from point-view, bird’s-eye view, and range-view representations. To improve robustness, it introduces a transformation consistency loss based on test-time augmentation to ensure agreement between original and augmented point clouds. PCSC-Net [119] combines point convolution and 3D sparse convolution for semantic segmentation. It generates large-size voxels from input point clouds, applies point convolution to extract voxel features, and then utilizes 3D sparse convolution to propagate features into neighbouring regions, enhancing feature extraction and context understanding. However, volumetric methods may lose information with low-resolution 3D grids, and their computational costs and memory requirements increase cubically with voxel resolution.

4.1.2 Direct Point-Based Methods

These methods operate directly on unstructured and irregular PCs, which poses a challenge for applying standard CNNs. PointNet [35] is a pioneering work in this domain, introducing a framework for processing direct PCs. Building upon PointNet, various approaches have been proposed, including pointwise MLP methods, point convolution methods, and graph-based methods, all aiming to enhance the processing and understanding of unstructured and unordered PC data.

Pointwise MLP methods: These methods utilize shared MLPs as the fundamental building block in their networks. However, the features extracted on a pointwise basis by these shared MLPs may face challenges in capturing the complex local geometry within PCs and the mutual interactions between points. To address these limitations, novel strategies have been introduced, including neighbouring feature pooling, attention-based aggregation, and local-global feature concatenation.

Neighbouring feature pooling: These methods are designed to capture local geometric patterns by aggregating information from nearby points to learn features for individual points in a PC. In [120], PointNet [35] was employed for semantic segmentation of elements such as pipes, valves, and background in several underwater environments. Additionally, the author created a novel PC dataset containing pipes and valve elements in various underwater scenarios. PointNet++ [52] performs a hierarchical grouping of points to learn features from large local regions. Subsequent developments, such as multi-scale grouping and multi-resolution grouping, have been introduced to address challenges arising from the non-uniform density of PCs. [121] proposed the surface-normal enhanced PointNet++ (SNEPointNet++) for semantic segmentation of defects, such as cracks and spalls, on concrete bridge surfaces. This approach emphasizes the utilization of normal vector, colour, and depth characteristics to address challenges associated with small size and imbalanced PC data. In [122], the authors introduced a focal loss function and a network named PC registration network (PCCR-Net), based on PointNet++, for segmenting precast concrete structures. PCCR-Net is specialized in segmenting various components such as columns, beams, slabs, walls, concrete, and rebars. Notably, the conventional negative log-likelihood loss function of PointNet++ was replaced with the focal loss function for gradient descent. Furthermore, the authors presented a synthetic PC dataset comprising diverse precast concrete components. [123] introduced ResPointNet++ featuring two NNs: a local aggregation operator for learning complex local structures and residual bottleneck modules to overcome gradient vanishing issues. ResPointNet++ demonstrates superior segmentation performance for indoor industrial systems compared to PointNet++, achieving F1-scores of 0.9874 and 0.6546, respectively. The PointSift module, as proposed by [124], achieves multi-scale representation by stacking and convolving features from the nearest points across eight different spatial orientations. This versatile module can seamlessly integrate into any PointNet-based framework, enhancing the network’s representation capability. [125] defines points neighbourhood in both the world space and feature space using K-means clustering and k-nearest neighbours (kNN) respectively. The learned point feature space is then structured by using pairwise distance loss and centroid loss. The mutual interaction between different points in the PC was explored by PointWeb [59] by constructing a local fully-linked web. An adaptive feature adjustment module is proposed to exchange information and refine features, followed by aggregating the learned features to obtain discriminative feature representation. RandLA-Net [126] introduces a lightweight NN designed to directly infer per-point semantics for large-scale PC segmentation tasks. The author incorporates a local feature aggregation module along with random point sampling to retain fine-grained geometric details during object segmentation. Multiple view pointwise networks (MVP-Net) [127] introduced space-filling curves and multi-rotation PC methods to expand the receptive field and aggregate the captured semantic features efficiently. Compared to RandLA-Net [126], MVP-Net demonstrates 11 times faster performance and higher efficiency in semantic segmentation tasks using the SemanticKITTI dataset. In another study, [128] investigated the impact of neighbourhood size selection on defect segmentation in 3D bridge PCs. The authors compare various sub-sampling approaches, including fast-graph, uniform, and random methods, to identify the optimal neighbourhood selection strategy.

Attention-based methods: These methods introduce innovative techniques for learning relations between points in PC data. [58] proposed a self-attention operator called GSA to learn relations between points. Later, the author used a task-agnostic sampling operation named GSS to replace the traditional FPS approach. This module is less sensitive

to outliers allowing a selective representative subset of points. The spatial distribution of the PC can be captured effectively by using the local spatial awareness network (LSA-Net) [129]. The LSA layer hierarchically generates spatial distribution weights based on relationships in spatial regions and local structures in the PC. Based on the CRF framework proposed by [103], [130] introduced an attention-based score refinement (ASR) module. This module computes weights for each point in the PC based on their initial segmentation scores, facilitating a refinement process where the scores of each point, along with those of its neighbours, are pooled together. The pooling operation is influenced by the computed weights, offering adaptability to efficiently integrate the module into various network architectures, thereby enhancing PC segmentation. [131] used an attention-based learning module for capturing local features and semantic relations in an anisotropic manner. Subsequently, a multi-scale context-guided aggregation module was used to differentiate points in the feature space, enhancing the scene-level understanding of semantic segmentation. Global attention network (GA-Net) [132] incorporated both point-independent and point-dependent GA modules for learning global contextual information across the entire PCs. Additionally, a point-adaptive aggregation block was introduced to group learned features, enhancing discriminative feature aggregation compared to linear skip connections. In [133], a semi-supervised learning (SmSL) approach called SPC-Net is introduced for segmenting various elements in tunnel PC data, including cables, segments, pipes, power tracks, supports, and tracks. Step-wise PC completion network (SPC-Net) utilizes a supervised learning model with attention mechanisms and a downsampling-upsampling structure to facilitate efficient learning and feature extraction. Furthermore, a formulated loss function is implemented to enable SPC-Net to conduct SmSL for multi-class object semantic segmentation of 3D tunnel PCs. [134] proposed a similar attention-based network called attention-enhanced sampling PC network (ASPC-Net), aimed at distinguishing defect classes in tunnel PC data. ASPC-Net incorporates a weighted focal loss strategy to overcome the impact of imbalanced data, enhancing its ability to accurately classify defects in tunnel PC datasets.

Local-global feature concatenation: This approach addresses the segmentation challenges posed by various object sizes and scales in large-scale PCs by integrating both local and global features. Many existing methods prioritize either global or local features, while hierarchical approaches often emphasize local features at the expense of global shape features. By concatenating local and global features, this approach enables comprehensive feature representation, enhancing segmentation accuracy across different object sizes and scales in large-scale PC datasets. PointMixer, as introduced in [135], facilitates information sharing among unstructured 3D PCs by substituting token-mixing MLPs with a SoftMax function. This method aggregates features across multiple points, encompassing intra-set, inter-set, and hierarchy sets, thereby promoting effective feature fusion and information exchange within the PC data. In [136], PS2-Net was introduced as a permutation-invariant approach for 3D semantic segmentation, integrating local structures and global context. The method leverages Edgeconv [8] to capture local structures and NetVLAD [137] to model global context from PCs, enabling comprehensive feature extraction for accurate segmentation. SCF-Net [138] presented a unique approach to learn spatial contextual features (SCF) tailored for large-scale PCs. SCF-Net uses a local polar representation (LPR) block to construct a representation invariant to z-axis rotation. Neighbouring representations are then aggregated via a dual-distance attentive pooling (DDAP) block to capture local features effectively. Furthermore, a global contextual feature (GCF) block utilizes both local and neighbourhood information to learn global context. SCF-Net's versatility allows it to be seamlessly integrated into encoder-decoder architectures for 3D semantic segmentation. LGFF-Net [139] introduces a novel local feature aggregation (LFA) module to capture geometric and semantic information concurrently, preserving original data integrity during cross-augmentation. Following this, a global feature extraction (GFE) module is used to extract global features. Ultimately, local and global features are concatenated using a U-shaped segmentation structure, enhancing overall segmentation performance. In [140], a dual feature complementary (DFC) module is proposed to learn local features effectively. This module employs a position-aware block to adaptively move with smaller point sets, enhancing the capture of geometric features. Additionally, a global correlation mining (GCM) module is utilized to gather contextual features, further improving semantic segmentation performance. In [141], the segmentation of overhead catenary systems in high-speed rails is enhanced by integrating local feature extraction with contextual feature information. Local features are extracted from both the local points and their neighbourhoods, followed by the aggregation of contextual information using CNN layers. Subsequently, feature enhancement and fusion techniques are applied to refine the segmentation process. Point central transformer (PointCT) [142] introduces a central-based attention mechanism and transformer architecture to address sparse annotations in PC semantic segmentation. Spatial positional encoding is introduced to focus on various geometries and scales for point representations, enhancing flexibility and enriching the representation of unlabeled points in PCs. In [3], a dempster-shafer (D-S) evidence-based feature fusion model was employed to integrate local and global features extracted from different CNN models. The study targeted the segmentation of tunnel PC defects, encompassing cable, pipe, segment, track, and power tracks. Results showcased enhanced segmentation scores compared to raw point-based segmentation models across various baselines. [143] introduced a transformer-based feature embedding network (3D Trans-Embed) for detecting defective industrial products. The method leverages a transformer model for PC segmentation and integrates local feature embedding technology along with multi-channel feature map fusion to enhance attention towards defective regions, thereby improving semantic segmentation outcomes.

Convolution-based Methods: These methods harness the intrinsic capabilities of CNNs to extract high-level discriminative features from complex spatial structures present in PCs. [144] used two CNNs and an RNN to conduct semantic segmentation of structural, architectural, and mechanical objects. The approach was trained and evaluated on PC data from 83 rooms, representing real-world industrial and commercial buildings. In the work by [145], a novel combinational convolutional block (CCB) called PCNet++ is introduced and applied to the synthetic gear dataset (Gear-PCNet++) to detect gear defects (wear, fracture, glue, and pitting) in manufacturing industries. PCNet++ replaces the convolution layer in MLP networks with the novel CCB to effectively extract local gear information while identifying its complex topology. The method outperforms PointNet, PointNet++, PointCNN, and KPConv on the gear PC dataset, achieving superior segmentation results. In [83], features of individual points within the PC are learned using a point-wise CNN for semantic segmentation and object recognition. Furthermore, parametric-continuous CNN (PCNN) [146] operates on non-grid data structures by employing a parameterized kernel function that spans continuous vector space. These methods utilize the spatial properties of PCs to develop point-based CNNs with spatial kernels, enabling the application of convolution operators tailored to the local structures of the PC. KPConv [80] presents a distinctive approach to 3D semantic segmentation by using radius neighbourhoods as input for convolution, ensuring a consistent receptive field. This method processes these neighbourhoods with weights determined spatially by a small set of kernel points. Additionally, KPConv incorporates a deformable operator to learn local shifts, enabling the customization of convolution kernels for improved alignment with the PC geometry. Dense connection-based kernel point network (DenseKP-NET) [147] extends the receptive field by introducing a multi-scale convolution kernel point module, facilitating the extraction of coarse-to-fine geometric features. Subsequently, a dense connection module refines these features while capturing contextual information. However, while kernel-based approaches excel in semantic segmentation, they may lack in providing ample local contextual features. To address this, [148] introduces attention kernel convolution (AKConv) to discern local contextual features while preserving object geometric shape information. In [149], a dilated point CNN (DP-CNN) is proposed to investigate the impact of the receptive field on existing point-convolution methods. DP-CNN enhances the receptive field size by aggregating features from dilated neighbours instead of KNN. Meanwhile, [150] introduces PointConvFormer, amalgamating point convolution with transformers to bolster model robustness. PointConvFormer utilizes pointwise CNN for feature extraction and computes attention weights based on feature disparities, refining convolutional weights and enhancing model performance.

Graph-based methods: These methods utilize a graph as the fundamental structure for applying convolution to irregular PCs. This approach eliminates the necessity of transforming PCs into regular grids or voxels, enabling direct processing on the intrinsic graph-like nature of PCs. [1] segmented the defects such as dents, protrusions or scratches on aircraft by using a region-growing network. The process involved initially smoothing the collected PC using a moving least squares (MLS) algorithm. Subsequently, curvature and normal information were collected for each point in the PC before applying the region-growing segmentation. DGCNN [8] treats neighbouring points as a local graph and feeds it into a filter-generating network to assign edge labels. Being a transformation invariance network, DGCNN is not affected by the order of local points. However, while it handles local points, it doesn't fully exploit the geometric information of neighbouring points in the PC. In [151], DGCNN was examined for segmenting concrete surface defects, where modifications to the loss function and data augmentation techniques, particularly flipping, were discussed to enhance performance. Later the author [9] proposed an improved DGCNN method for defect segmentation on concrete surfaces, leveraging normal vectors and depths to detect surface defects (cracks, and spalls) effectively. In the study conducted by [152], the segmentation of bridge components (abutment, girder, background, pier, deck, slab, and surface) for inspection was performed using PointNet, PointCNN, and DGCNN. The results showed that DGCNN outperformed other networks, achieving an OA and mIoU of 0.94 and 0.86, respectively. Unlike PointNet, which focuses solely on global features of input points, DGCNN incorporates information from neighbouring points, enabling the generation of meaningful features to classify different bridge components based on their relationships with the surroundings. In [153] presented an enhanced version of DGCNN incorporating additional features such as normals and colours, to segment architectural elements (arc, column, decoration, floor, door, wall, window, stairs, vault, and roof) in buildings. The study demonstrated superior segmentation performance of the enhanced DGCNN compared to PointNet, PointNet++, PCNN, and the original DGCNN, respectively. In [10], a graph attention convolution (GAC) with learnable kernels was introduced, enabling dynamic adaptation to the structure of objects. GAC effectively learns discriminative features for semantic segmentation, offering similar characteristics to traditional CRF models. [154] extended the concept of GAC with the introduction of a cross-scale graph interaction network (CGIN) for segmenting remote-sensing images. CGIN used a CGI module to extract multi-scale semantic features and a boundary feature extraction (MBFE) module to learn multi-scale boundary features. Furthermore, a similarity-guided aggregation module calculates the similarity between these features, highlighting boundary information within semantic features. In [155], a simulation-to-real transfer learning (TL) approach is introduced, utilizing DGCNN as the backbone network for segmenting industrial elements such as pole pot, electric connection, gear container, cover, screws, magnets, armature, lower gear, and upper gear. The author also introduced a patch-based attention network to tackle imbalanced learning challenges. [156] introduced a local-attention edge convolution (LEA-Conv) layer to construct a local graph by considering neighbourhood

points along sixteen directions. The LAE-Conv layer assigns attention coefficients to each edge of the graph while aggregating the extracted point features through a weighted sum computation of its neighbourhood. This local attention mechanism effectively captures long-range spatial contextual features, thereby enhancing the precision of semantic segmentation. [157] proposed a local-global graph CNN for semantic segmentation to capture both short and long-range dependencies within PCs. The author computes a weighted adjacency matrix for the local graph, utilizing information from neighbouring points, and performs feature aggregation to capture spatial geometric features. Subsequently, these learned features are fed into a global spatial attention module to extract long-range contextual information.

Table 3: Summary of PC-based defect segmentation in industrial systems.

Ref.	Application	Classes	Method	Results	Points/Objects
[153]	Semantic segmentation of architectural elements	10 classes: column, decoration, door, arc, wall, window, stairs, vault and roof	PointNet	OA = 21.6, mIoU = 10.9	114 M
			PointNet++	OA = 24.5, mIoU = 18.0	
			DGCNN	OA = 54.7, mIoU = 35.8	
			PCNN	OA = 39.5, mIoU = 33.1	
			Modified DGCNN	OA = 71.6, mIoU = 37.7	
[152]	Semantic segmentation of architectural elements	6 classes: abutment, slab, pier, girder, surface, and background	PointNet	OA = 93.8, mIoU = 84.3	N/A
			PointCNN	OA = 92.6, mIoU = 76.8	
			DGCNN	OA = 94.5, mIoU = 86.9	
[158]	Segmentation of bridge elements	3 classes: deck, pier, and background	PointNet	OA = 94, mIoU = 84	50,000
[120]	Semantic segmentation of underwater pipe	3 classes: pipe, valve, and background	PointNet	F1 score = 89.3	262
[151]	Segmentation of concrete surface	3 classes: crack, spall, and normal	DGCNN	OA = 98, F1 score = 98	49 M
[110]	Segmentation of synthetic concrete defects	2 classes: cracks, and spalls	InspectionNet	mAcc = 96.46	12,000
[95]	Segmentation of bridge elements	2 classes: slab and pier	DGCNN	OA = 95.9, mIoU = 71.1	447 M
			PointNet	OA = 84.4, mIoU = 45.9	
[121]	Semantic segmentation of concrete bridge elements	3 classes: cracks, spalls, and normal	SNEPointNet++	OA = 95.9, mIoU = 83.26	27 M
			Adaptive PointNet++	OA = 97.12, mIoU = 63.36	
[145]	Segmentation of gear	5 classes: basic gear, fracture, glue, wear, pitting	Gear-PCNet++	OA = 99.53, mIoU = 98.97	10,000
			PointNet++	OA = 99.29, mIoU = 98.50	
			PointCNN	OA = 99.43, mIoU = 98.76	
			KPCConv	OA = 99.64, mIoU = 97.50	
[122]	Semantic Segmentation in precast concrete rebar	4 classes: column, beam, slab, and wall	PCCR-Net	OA = 97.47, mIoU = 93.12	342
			PointNet++	OA = 95.17, mIoU = 87.68	
[159]	Panoptic segmentation in railway infrastructure	7 classes: informative signs, masts, traffic lights, traffic signs, cables, droppers and rails	PointNet++	OA = 95.34, mIoU = 80.3	4.5 M
[133]	Semi-supervised segmentation of 3D tunnel elements	6 classes: cable, segment, pipe, power track, support, track	SPCNet	OA = 97.23, mIoU = 97.41	32 M
[134]	Segmentation of 3D tunnel elements	7 classes: cable, segment, pipe, power track, seepage, support, track	ASPCNet	OA = 97.58, mIoU = 89.80	34 M
[160]	Segmentation of 3D tunnel elements	7 classes: cable, segment, pipe, power track, seepage, support, track	DGCNN	F1 = 91.9, mIoU = 97.5	34 M
[141]	Segmentation of overhead catenary systems (OCS's) in high-speed rails	8 classes: cantilevers, catenary wires, contact wires, droppers, insulators, poles, registration arms, steady arms	PointNet KNN+CNN	F1 = 98.1, mIoU = 96.3 Precision = 97.50, mIoU = 94.84	16 M
			PointNext	Precision = 96.49, mIoU = 93.39	
			PointNet++ (SSG)	Precision = 96.42, mIoU = 93.06	
[155]	Segmentation of industrial elements	9 classes: pole pot, electric connection, gear container, cover, screws, magnets, armature, lower gear, upper gear	PointNet DGCNN	Precision = 94.52, mIoU = 89.18 mIoU (real) = 93.75, mIoU (Simulation) = 98.01	5.2 M
[123]	Industrial indoor LiDAR dataset	6 classes: I-beam, pipe, pump, rectangular beam, and tank	PointNet	OA = 53.0, mIoU = 21.1	5 M
			PointNet++	OA = 70.6, mIoU = 45.5	
[143]	Segmentation of vegetation and industrial products	6 classes: potatoes, carrots, peaches, cookies, bagels, cable, 6 industrial products: cable gland, dowel, tyres, foam, and ropes	ResPointNet++ 3D Trans-Embed	OA = 94.0, mIoU = 87.3 F1-score = 83.32, Precision = 87.82	4,000
			PCT	F1-score = 72.71, Precision = 73.48	
			PointNet++	F1-score = 65.15, Precision = 69.18	
[161]	Instance segmentation of different object shapes in oil refinery, petrochemical plant and warehouse	8 classes: cylinders, angles, channels, I-beams, elbows, flanges, valves and miscellaneous	CLOI-NET	mPrec = 73.2, mRec = 71.1	N/A
[144]	Semantic segmentation of structural, architectural, and mechanical objects	6 classes: beam, ceiling, column, floor, pipe, and wall	ASIS CNN+RNN	mPrec = 74, mRec = 24.9 mAcc = 86.13	10.8 M
[162]	Segmentation of surface defects in fibre composites	2 classes: void defects and surface defects	CNN MaskPoint	mAcc = 84.70 mAcc = 99.97, mIoU = 94.02	120 M
			PointNet++	mAcc = 99.50, mIoU = 58.81	
			PointNet	mAcc = 99.41, mIoU = 62.72	
			KPCConv	mAcc = 99.32, mIoU = 49.53	
[64]	Defect segmentation in concrete sewer pipes	5 classes: 3 circular defects of varying diameter, square and triangular defect	PointTransformer Improved PointNet++ PointNet++ Point Transformer	mAcc = 99.38, mIoU = 49.83 mIoU = 94.15 mIoU = 82.69 mIoU = 86.31	1.4 M

4.2 Instance Segmentation

Instance segmentation, in contrast to semantic segmentation, presents a more challenging task as it requires distinguishing points sharing the same semantic meaning. To address this complexity, instance segmentation methods fall into

two main categories: proposal-based methods and proposal-free methods. Figure 8 provides a comparison between proposal-based and proposal-free instance segmentation methods using the 3D ScanNet dataset [18].

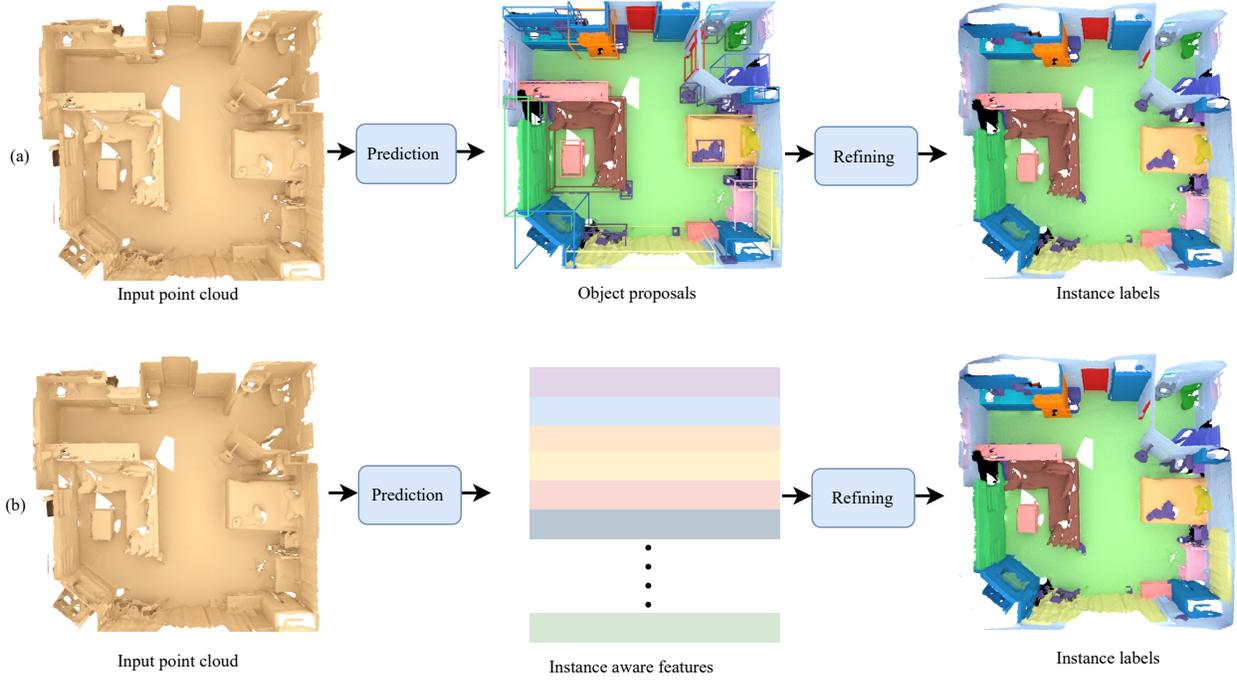


Figure 8: Illustration of 3D instance segmentation frameworks on ScanNet benchmark [18]: (a) Proposal-based methods, and. (b) proposal-free methods. Recreated from [18]

4.2.1 Proposal-Based Methods

Proposal-based instance segmentation methods can be conceptualized as a fusion of object detection and mask prediction strategies. These methods follow a top-down pipeline where the initial step involves the generation of region proposals, usually bounding boxes (BBox's) followed by predicting instance masks within these proposed regions. The pipeline encompasses multiple stages, including proposal generation, classification, and mask prediction, often integrating techniques from both object detection and semantic segmentation. In [162], the authors propose Mask-Point, a multi-head region proposal extractor to generate multiple regions of interest, allowing networks to focus on potential defective regions. Following this, an aggregation module is designed to improve the segmentation of surface defects in fibre-reinforced composites. [158] proposed a region-CNN (R-CNN) method by combining region proposals with features extracted from CNN to segment cracks on concrete bridges. 3D-SIS [163] is an FCN designed for 3D semantic instance segmentation using RGB-D scans. This network uses a series of CNN layers to extract 2D features for each pixel, which are then projected back onto 3D voxel grids. The RGB-D scan features are processed by 3D-CNN and aggregated into a global semantic feature map. Subsequently, 3D-Region Proposal Network (3D-RPN) and 3D Region of Interest (3D-ROI) layers are utilized to predict the locations of BBox's, instance masks, and object class labels. Building upon 3D-SIS, [164] applies this framework to segment casting defected regions (CDR) in a foundry industrial plant, introducing a non-linear topological dimension parameter to characterize the geometrical features of the segmented regions. Generative shape proposal network (GSPN) [165] introduces a novel approach for proposal generation by reconstructing shapes from scenes, contrasting with conventional methods that regress BBox's. These generated proposals undergo refinement through a region-based PointNet, with the final labels determined by predicting point-wise binary masks for each class label. Importantly, GSPN incorporates a mechanism to discard trivial proposals by directly learning geometric features from the PCs. Based on PointNet++, [166] introduced 3D-BoNet, a single-stage, anchor-free, and end-to-end trainable method for achieving instance segmentation on PCs. 3D-BoNet adopts a direct regression approach to predict 3D BBox's for all instances in a PC, while simultaneously predicting point-level masks for each instance. Gaussian instance center network (GICN) [167] utilizes Gaussian heat maps to represent the locations of instance centres distributed across the scene. By estimating the size of each instance, GICN adjusts its feature extraction process to capture relevant information within the specified neighbourhood, thereby enhancing the precision and adaptability of segmentation. In [168], OccuSeg, an occupancy-aware 3D instance segmentation method, was

introduced to predict point-wise instance-level segmentation. It leverages a 3D occupancy signal to predict the number of occupied pixels/voxels for each instance. This occupancy signal, learned in conjunction with feature and spatial embeddings, guides the clustering stage of 3D instance segmentation, enhancing segmentation accuracy.

Transformer-Based Methods: Transformer-based methods have emerged as powerful tools in various CV tasks, including semantic segmentation and instance segmentation of PCs. These methods utilize the self-attention mechanism to capture long-range semantic relationships within the input data, combining both positional and feature information effectively. [169] utilizes a transformer architecture to compute object features directly from the PC data while refining predictions by updating the spatial encoding of the objects across different stages. On the other hand, segmenting objects with transformers (SOTR) [170] combines the strengths of both CNN and transformer methodologies for segmenting objects. This is achieved by using a feature pyramid network (FPN) alongside twin transformers to extract lower-level features and capture long-range context dependencies for object segmentation. In the medical domain, [171] introduced a fusion of CNN and transformers termed CoTr, for 3D-multi-organ segmentation. In this approach, CNNs were used for feature extraction, while a deformable transformer was utilized to capture long-range dependencies within high-resolution and multi-scale feature maps, enhancing the segmentation performance. BoundaryFormer [172] used pixel-wise masks as ground truth to predict object boundaries in the form of polygons. The method evaluates the loss using an end-to-end differentiable rasterization model, enabling precise delineation of object boundaries during instance segmentation. SPFormer, as proposed by [173], is an end-to-end two-stage method designed for instance segmentation of PCs. In the first stage, potential features extracted from the input PCs are aggregated into super points. Subsequently, a query decoder equipped with transformers is used to directly predict instances based on these super points, facilitating efficient and accurate instance segmentation. Mask3D [174] used stacked transformer decoders to predict instance queries, enabling the encoding of both semantic and geometric information for individual instances within a scene. While previous methods relied on instance masks for computing object queries followed by iterative refining, which often led to slow convergence, Mask3D offers an alternative approach. To alleviate the dependency on mask attention, [175] proposed a mask-attention-free transformer (MAFTr). MAFTr utilizes contextual relative position encoding for cross-attention, where position queries are iteratively updated to provide more accurate representations, enhancing the efficiency and effectiveness of instance segmentation.

Indeed, proposal-based methods for instance segmentation offer an intuitive approach by combining object detection and mask prediction strategies. However, these methods typically involve multi-stage training processes and the need for pruning redundant proposals, which can be time-consuming and computationally expensive. This complexity arises from the necessity of generating region proposals, such as BBoxes, followed by classification and mask prediction within these proposed regions. As a result, while proposal-based methods may achieve high accuracy, they often come with a significant computational cost and training overhead.

4.2.2 Proposal-Free Methods

Proposal-free methods for instance segmentation leverage the inherent characteristics of point clouds, such as their spatial distribution and semantic information, without relying on explicit region proposals. Instead, these methods typically use clustering techniques to group points with similar semantic meanings into distinct instances. By directly segmenting PCs into instances without the need for explicit proposals, proposal-free methods can be more computationally efficient and simpler in concept compared to proposal-based approaches. Similarity group proposal network (SGPN) [176] is a pioneering work designed to learn features and semantic maps for individual points in a PC. This network constructs a similarity matrix which encapsulates the similarity between every pair of features within the PC. To enhance the discriminative features, SGPN uses a double-hinge loss, which adjusts both the similarity matrix and the semantic segmentation results. Later, it uses a heuristic non-maximal suppression technique to merge similar points into distinct instances. However, the construction of the similarity matrix demands substantial memory resources, limiting the scalability of this method. Multi-scale affinity with sparse convolution (MASC) [177] utilizes sparse convolution to predict semantic scores for each voxel while capturing the point affinity between neighbouring voxels across multiple scales. Furthermore, it uses a clustering algorithm to organize points according to the learned local similarities and the inherent mesh topology. [178] proposed a structure-aware loss function to learn discriminative embeddings for each instance by considering the similarity between both geometric and embedding information. The author proposed attention-based kNN to refine the learned features by grouping information from neighbours while eliminating the quantization error caused by the 3D voxel.

Several methods have been proposed by integrating semantic category and instance label prediction into a single task. Milestones in 3D PC instance segmentation, including both proposal-based and proposal-free methods, are depicted in Figure 9. [179] integrates the advantages of both instance and semantic segmentation through an end-to-end learnable module called associatively segmenting instances and semantics (ASIS). The ASIS module incorporates semantic-aware point-level embedding to achieve instance segmentation and performs instance fusion to obtain semantic segmentation simultaneously. [180] introduced a joint instance and semantic segmentation (JISS) module, which combines both

instance and semantic segmentation to generate discriminative features. To address the large memory consumption of JSNet, the authors proposed dynamic filters for convolution (DFConv) on PCs. Based on JSNet, DFConv, and an enhanced JISS (JISS*) module, [181] introduced JSNet++ to enhance instance segmentation. 3D-Multi proposal aggregation (3D-MPA) [182] presents a technique for predicting object proposals by utilizing semantic features derived from a sparse volumetric backbone network. In contrast to conventional non-maximum suppression (NMS), this method employs the MPA strategy, based on learned features, to derive semantic instances from the generated object proposals.

Grouping-Based Methods: In contrast to proposal-based methods, grouping-based methods follow a bottom-up pipeline approach. These methods learn point-wise semantic labels and instance centre offsets. Subsequently, the offset points and semantic predictions are aggregated to form instances [183, 184]. [185] proposed a multi-task algorithm (MSA) to learn unique feature embeddings for each instance by leveraging grouping or clustering information associated with individual objects. PointGroup [186] focuses on grouping points by identifying the void space between distinct objects. To achieve this, the authors proposed a two-branch network capable of extracting point features, predicting semantic labels, and computing offsets concurrently. These offsets are then employed to relocate each point towards its corresponding instance centroid. Based on PointGroup, [187] introduced a clustering-based framework called hierarchical aggregation for 3D IS (HAIS) to produce detailed instance predictions while also effectively filtering out noisy points within instance predictions. Dyco3D [188] introduced dynamic convolution kernels, which encode category-specific context by utilizing a sub-network to explore homogenous points showing close votes for instance centroids and sharing the semantic labels. The parallel decoding of instance masks is accomplished by convolving the generated class-specific filters with coordinate information. SST-Net [183] introduced a semantic super-point tree, where each super-point represents a geometrically homogeneous neighbourhood. This method utilizes tree traversal for object proposal by splitting non-similar nodes in this semantic super-point tree. SoftGroup [184] addresses errors arising from hard semantic predictions by performing grouping based on semantic scores. The method uses a top-down refinement module using U-Net to improve positive samples while suppressing false positives introduced by incorrect semantic predictions.

In summary, while proposal-free methods alleviate the computational burden associated with region-proposal mechanisms, they often exhibit lower objectness in the resulting grouped instance segments. This limitation stems from their inherent inability to explicitly detect object boundaries, leading to less precise delineation of individual objects within the point cloud.

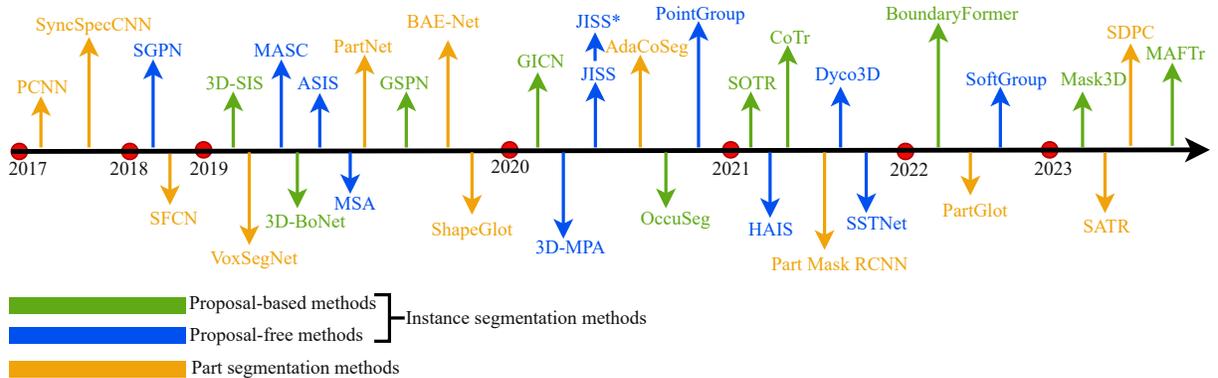


Figure 9: Chronological overview of the most relevant DL-based 3D instance and part segmentation methods.

4.3 Part Segmentation

Part segmentation involves categorizing the PC into distinct groups, where each group represents a specific physical part of the object. However, part segmentation encounters two significant challenges. First, parts with the same semantic label may exhibit considerable geometric variation and ambiguity. Second, objects with identical semantic meanings may consist of different numbers of constituent parts. Several milestone 3D PC part segmentation methods have been illustrated in Figure 9. VoxSegNet [189] introduced a spatial dense extraction (SDE) module to extract multi-scale discriminative features from sparse volumetric data. These learned features are contextually selected and aggregated through an attention feature aggregation (AFA) module, ensuring dense prediction with semantic consistency and enhanced accuracy. PartNet [190] introduces a top-down, fine-grained, and hierarchical approach to part segmentation. Unlike conventional methods that segment shapes into a fixed set of labels, PartNet formulates part segmentation as a cascade binary labeling process. This methodology decomposes the input PC into an arbitrary number of parts,

determined by the underlying geometric structures. [191] introduced an end-to-end network called projective CNNs (PCNNs), which combines FCNs and surface-based CRFs to achieve part segmentation of 3D shapes. The authors selected images from multiple views to ensure optimal surface coverage and fed them into the network to generate per-part confidence maps. These confidence maps are then aggregated using surface-based CRFs to label the entire surface. However, dealing with different shapes resulted in different nearest-neighbour graphs in the PC, posing challenges for weight sharing among convolution kernels across various shapes. To address this challenge, synchronized spectral CNN (SyncSpecCNN) [192] uses a spectral network for convolution, allowing weight sharing across different non-isometric shapes. Additionally, [193] introduced part segmentation on 3D meshes using shape FCNs (SFCNs). The author utilized SFCNs to process low-level geometric features and refined the segmentation outcomes through feature voting-based multi-level graph cuts. In [194], the authors proposed Part-Mask RCNN for predicting shape categories, bounding boxes, object masks, and object part masks in RGB-D images. The authors utilized a voting-based pose estimation algorithm on semantic information of the objects to obtain part segmentation. [195] proposed an adaptive shape co-segmentation (AdaCoSeg) network to address the challenges associated with re-training and adapting to newer input datasets. AdaCoSeg takes a set of unsegmented PC shapes as input and iteratively minimizes the group consistency loss to produce shape part labels. Unlike traditional CFR methods, the authors refine and denoise the part proposals using a pre-trained part-refinement network. The branched auto-encoder network (BAE-NET) [196] tackles the 3D shape co-segmentation task by framing it as a representation learning challenge. The goal of this approach is to discover the most concise part representations by minimizing the shape reconstruction loss. Using an encoder-decoder architecture, each branch of the network is dedicated to learning a condensed representation for a particular part shape. The features acquired from each branch, combined with the point coordinates, are fed into the decoder to produce a binary value indicating whether the point belongs to that part. In the medical domain, [197] proposed a shape-aware segmentation (SAS) technique for processing MRI imaging scans. This method imposes geometric constraints on both labeled and unlabeled input data. It involves learning a shape-aware representation using a signed distance map (SDM) approach. Following this, the obtained predictions undergo refinement through an adversarial loss. Based on the ShapeGlot framework [198], PartGlot [199] utilizes a transformer-based attention mechanism to understand the regions corresponding to semantic parts by leveraging linguistic descriptions. [200] proposed a soft density peak clustering (SDPC) algorithm tailored for 3D shape segmentation. [201] developed a segmentation assignment with topological re-weighting (SATR) to achieve part segmentation from the predicted multi-view BBox's. Firstly, gaussian geodesic re-weighting is performed to adjust weights by considering the geodesic distance from potential segment centres. Secondly, a graph kernel is used to refine the inferred weights considering the neighbour's visibility. These two techniques are combined to achieve state-of-the-art 3D shape segmentation for fine-grain queries. [201] used geodesic curves for discriminative modeling of the object shapes within an NN framework. The method involves selecting pairs of 3D points on depth images to compute surface geodesics. The approach leverages a large training set of geodesics created using minimal ground truth instance annotations, where each geodesic is labeled binary to indicate whether it belongs entirely to one instance segment. An NN is then trained to classify geodesics based on these labels. During inference, geodesics are generated from selected seed points in the test depth image, and a convex hull is constructed for points classified by the neural network (NN) as belonging to the same instance, thereby achieving instance segmentation.

In summary, 3D part segmentation allows a fine-grained understanding of objects by categorizing the PCs into distinct parts, providing geometric information about the objects in the scene. Also, it contributes to semantic interpretation, enabling systems to recognize and label structural components of objects. However, objects from the same semantic class may have significant intra-class variability, making it challenging to define consistent part boundaries. Tables 3 present the outcomes of defect segmentation for industrial systems.

4.4 Summary

This section presents key challenges and research directions in processing 3D PCs for defect classification and segmentation:

- The survey suggests that projection-based methods often adopt network architectures similar to their 2D image counterparts. However, a key limitation of these methods is the loss of information due to the conversion from 3D to 2D projection. On the other hand, volumetric-based representations encounter challenges with significantly increased computational and memory costs, attributed to the cubic growth in resolution. Addressing these issues, sparse convolution methods leveraging indexing structures emerge as a promising solution that needs to be further explored.
- Popular for defect classification and segmentation, point-based networks lack explicit neighbouring information, relying on costly neighbour-searching mechanisms, thus limiting the efficiency.

- Learning from imbalanced data remains a challenge, with approaches struggling in minority classes despite strong overall performance. Novel techniques are needed to handle imbalanced datasets effectively, such as data augmentation, class balancing, or specialized loss functions.
- Domain adaptation and transfer learning can overcome the need for extensive labeled datasets, especially in scenarios with limited labeled data availability. Techniques like generative adversarial networks (GANs) or other generative models can augment datasets, expanding training data and improving model generalization.
- The literature presents numerous studies dedicated to defect segmentation within general objects or space using 3D PC data. However, a significant gap exists in the research concerning the detection of damages within industrial systems using 3D PC data. Despite the promising outcomes of semantic segmentation in PC analysis for damage detection in industrial systems, its effectiveness heavily relies on the availability of comprehensive datasets for model training. Unfortunately, the literature lacks sufficient datasets for defect estimation in industrial systems, underscoring the urgent need to collect abundant and efficient data for this purpose.
- Instance segmentation represents a challenging task in computer vision, combining aspects of target detection and semantic segmentation. While limited studies are focusing on 3D instance segmentation of defects in industrial systems, the future holds promising prospects for the development of DL models in this domain.

While there has been considerable research in PC shape classification and object segmentation across fields like robotics, autonomous vehicles, and other computer vision applications using 3D PC data, the detection of damages in industrial systems remains relatively underexplored. This gap signifies a significant opportunity for future research to devise specialized methods and models specifically addressing the distinct challenges of condition monitoring in industrial systems.

5 Conclusion

This paper provides a comprehensive survey and discussion of DL-based methods for PC classification and segmentation in recent years. The introduction outlines the significance of PCs and their applications, highlighting the challenges associated with processing this unique form of data. The publicly available 3D PC datasets for object classification and segmentation have been thoroughly discussed. The review paper encompasses various DL-based approaches for PC defect shape classification, categorizing them into view-based methods, volumetric-based methods, and direct PC methods. The paper presents a comparison of the performance of these existing methods, providing insights into their strengths and limitations for industrial systems. Finally, the discussion section explores the prospects and potential research directions in this field, contributing to a holistic understanding of the advancements and challenges in DL-based PC analysis.

The PC classification and segmentation method finds applications in various real-world scenes, including indoor environments, roads, railways, buildings, etc. However, the diversity of these scenes poses challenges in determining the specific advantages of numerous PC classification methods. Consequently, researchers face the task of selecting a classification algorithm that aligns with the requirements of a given scenario, emphasizing the need for adaptability to real-world conditions. This challenge is further underscored by the scarcity of suitable datasets, as discussed earlier. The necessity for researchers to choose an appropriate classification algorithm based on the specific characteristics of the scene highlights the ongoing issue of dataset shortages in the field. This points towards the importance of expanding and diversifying datasets to better evaluate and improve the efficacy of PC classification methods across different real-world scenarios.

Declaration of competing interest

The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

Acknowledgment

This research was supported by Aalborg University, Liftra ApS (Liftra), and Dynamica Ropes ApS (Dynamica) in Denmark under the Energiteknologiske Udviklings- og Demonstrationsprogram (EUDP) program through project grant number 64021-2048.

References

- [1] Igor Jovančević, Huy-Hieu Pham, Jean-José Orteu, Rémi Gilblas, Jacques Harvent, Xavier Maurice, and Ludovic Brèthes. 3d point cloud analysis for detection and characterization of defects on airplane exterior surface. *Journal of Nondestructive Evaluation*, 36:1–17, 2017.
- [2] Yulan Guo, Hanyun Wang, Qingyong Hu, Hao Liu, Li Liu, and Mohammed Bennamoun. Deep learning for 3d point clouds: A survey. *IEEE transactions on pattern analysis and machine intelligence*, 43(12):4338–4364, 2020.
- [3] Huang Zhang, Changshuo Wang, Shengwei Tian, Baoli Lu, Liping Zhang, Xin Ning, and Xiao Bai. Deep learning-based 3d point cloud classification: A systematic survey and outlook. *Displays*, page 102456, 2023.
- [4] Qian Wang and Min-Koo Kim. Applications of 3d point cloud data in the construction industry: A fifteen-year review from 2004 to 2018. *Advanced Engineering Informatics*, 39:306–319, 2019.
- [5] Yifan Feng, Zizhao Zhang, Xibin Zhao, Rongrong Ji, and Yue Gao. Gvcnn: Group-view convolutional neural networks for 3d shape recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 264–272, 2018.
- [6] Peng-Shuai Wang, Yang Liu, Yu-Xiao Guo, Chun-Yu Sun, and Xin Tong. O-cnn: Octree-based convolutional neural networks for 3d shape analysis. *ACM Transactions On Graphics (TOG)*, 36(4):1–11, 2017.
- [7] Shaoshuai Shi, Xiaogang Wang, and Hongsheng Li. Pointcnn: 3d object proposal generation and detection from point cloud. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 770–779, 2019.
- [8] 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.
- [9] Fardin Bahreini and Amin Hammad. Dynamic graph cnn based semantic segmentation of concrete defects and as-inspected modeling. *Automation in Construction*, 159:105282, 2024.
- [10] Lei Wang, Yuchun Huang, Yaolin Hou, Shenman Zhang, and Jie Shan. Graph attention convolution for point cloud semantic segmentation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 10296–10305, 2019.
- [11] Shuo Chen, Tan Yu, and Ping Li. Mvt: Multi-view vision transformer for 3d object recognition. *arXiv preprint arXiv:2110.13083*, 2021.
- [12] Jie Li, Zhao Liu, Li Li, Junqin Lin, Jian Yao, and Jingmin Tu. Multi-view convolutional vision transformer for 3d object recognition. *Journal of Visual Communication and Image Representation*, 95:103906, 2023.
- [13] Jianhui Yu, Chaoyi Zhang, Heng Wang, Dingxin Zhang, Yang Song, Tiange Xiang, Dongnan Liu, and Weidong Cai. 3d medical point transformer: Introducing convolution to attention networks for medical point cloud analysis. *arXiv preprint arXiv:2112.04863*, 2021.
- [14] Daniel Munoz, J Andrew Bagnell, Nicolas Vandapel, and Martial Hebert. Contextual classification with functional max-margin markov networks. In *2009 IEEE Conference on Computer Vision and Pattern Recognition*, pages 975–982. IEEE, 2009.
- [15] Franz Rottensteiner, Gunho Sohn, Jaewook Jung, Markus Gerke, Caroline Baillard, Sebastien Benitez, and Uwe Breitkopf. The isprs benchmark on urban object classification and 3d building reconstruction. *ISPRS Annals of the Photogrammetry, Remote Sensing and Spatial Information Sciences; I-3*, 1(1):293–298, 2012.
- [16] Andrés Serna, Beatriz Marcotegui, François Goulette, and Jean-Emmanuel Deschaud. Paris-rue-madame database: a 3d mobile laser scanner dataset for benchmarking urban detection, segmentation and classification methods. In *4th international conference on pattern recognition, applications and methods ICPRAM 2014*, 2014.
- [17] Bruno Vallet, Mathieu Brédif, Andrés Serna, Beatriz Marcotegui, and Nicolas Paparoditis. Terramobilita/iqmulus urban point cloud analysis benchmark. *Computers & Graphics*, 49:126–133, 2015.
- [18] Angela Dai, Angel X. Chang, Manolis Savva, Maciej Halber, Thomas Funkhouser, and Matthias Nießner. Scannet: Richly-annotated 3d reconstructions of indoor scenes. In *Proc. Computer Vision and Pattern Recognition (CVPR)*, IEEE, 2017.
- [19] Iro Armeni, Ozan Sener, Amir R Zamir, Helen Jiang, Ioannis Brilakis, Martin Fischer, and Silvio Savarese. 3d semantic parsing of large-scale indoor spaces. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1534–1543, 2016.
- [20] Timo Hackel, Nikolay Savinov, Lubor Ladicky, Jan D Wegner, Konrad Schindler, and Marc Pollefeys. Semantic3d.net: A new large-scale point cloud classification benchmark. *arXiv preprint arXiv:1704.03847*, 2017.

- [21] Xavier Roynard, Jean-Emmanuel Deschaud, and François Goulette. Paris-lille-3d: A large and high-quality ground-truth urban point cloud dataset for automatic segmentation and classification. *The International Journal of Robotics Research*, 37(6):545–557, 2018. doi:10.1177/0278364918767506.
- [22] A. Geiger, P. Lenz, and R. Urtasun. Are we ready for Autonomous Driving? The KITTI Vision Benchmark Suite. In *Proc. of the IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, pages 3354–3361, 2012.
- [23] Weikai Tan, Nannan Qin, Lingfei Ma, Ying Li, Jing Du, Guorong Cai, Ke Yang, and Jonathan Li. Toronto-3D: A large-scale mobile lidar dataset for semantic segmentation of urban roadways. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops*, pages 202–203, 2020.
- [24] Nina Varney, Vijayan K Asari, and Quinn Graehling. Dales: A large-scale aerial lidar data set for semantic segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops*, pages 186–187, 2020.
- [25] Holger Caesar, Varun Bankiti, Alex H. Lang, Sourabh Vora, Venice Erin Liong, Qiang Xu, Anush Krishnan, Yu Pan, Giancarlo Baldan, and Oscar Beijbom. nuscenes: A multimodal dataset for autonomous driving, 2020.
- [26] Zhirong Wu, Shuran Song, Aditya Khosla, Fisher Yu, Linguang Zhang, Xiaoou Tang, and Jianxiong Xiao. 3d shapenets: A deep representation for volumetric shapes. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1912–1920, 2015.
- [27] Angel X Chang, Thomas Funkhouser, Leonidas Guibas, Pat Hanrahan, Qixing Huang, Zimo Li, Silvio Savarese, Manolis Savva, Shuran Song, Hao Su, et al. Shapenet: An information-rich 3d model repository. *arXiv preprint arXiv:1512.03012*, 2015.
- [28] Jiachen Sun, Qingzhao Zhang, Bhavya Kailkhura, Zhiding Yu, Chaowei Xiao, and Z. Morley Mao. Benchmarking robustness of 3d point cloud recognition against common corruptions. *arXiv preprint arXiv:2201.12296*, 2022.
- [29] Mikaela Angelina Uy, Quang-Hieu Pham, Binh-Son Hua, Duc Thanh Nguyen, and Sai-Kit Yeung. Revisiting point cloud classification: A new benchmark dataset and classification model on real-world data. In *International Conference on Computer Vision (ICCV)*, 2019.
- [30] Meida Chen, Qingyong Hu, Zifan Yu, Hugues THOMAS, Andrew Feng, Yu Hou, Kyle McCullough, Fengbo Ren, and Lucio Soibelman. Stpls3d: A large-scale synthetic and real aerial photogrammetry 3d point cloud dataset. In *33rd British Machine Vision Conference 2022, BMVC 2022, London, UK, November 21-24, 2022*. BMVA Press, 2022. URL <https://bmvc2022.mpi-inf.mpg.de/0429.pdf>.
- [31] Shuran Song, Samuel P Lichtenberg, and Jianxiong Xiao. Sun rgb-d: A rgb-d scene understanding benchmark suite. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 567–576, 2015.
- [32] Mike Roberts, Jason Ramapuram, Anurag Ranjan, Atulit Kumar, Miguel Angel Bautista, Nathan Paczan, Russ Webb, and Joshua M Susskind. Hypersim: A photorealistic synthetic dataset for holistic indoor scene understanding. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 10912–10922, 2021.
- [33] Hang Su, Subhransu Maji, Evangelos Kalogerakis, and Erik Learned-Miller. Multi-view convolutional neural networks for 3d shape recognition. In *Proceedings of the IEEE international conference on computer vision*, pages 945–953, 2015.
- [34] Bonsang Koo, Raekyu Jung, Youngsu Yu, and Inhan Kim. A geometric deep learning approach for checking element-to-entity mappings in infrastructure building information models. *Journal of Computational Design and Engineering*, 8(1):239–250, 2021.
- [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.
- [36] Zifan Shao, Kuangrong Hao, Bing Wei, and Xue-Song Tang. Solder joint defect detection based on depth image cnn for 3d shape classification. In *2021 CAA symposium on fault detection, supervision, and safety for technical processes (SAFEPROCESS)*, pages 1–6. IEEE, 2021.
- [37] Tan Yu, Jingjing Meng, and Junsong Yuan. Multi-view harmonized bilinear network for 3d object recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 186–194, 2018.
- [38] Chao Ma, Yulan Guo, Jungang Yang, and Wei An. Learning multi-view representation with lstm for 3-d shape recognition and retrieval. *IEEE Transactions on Multimedia*, 21(5):1169–1182, 2018.
- [39] Zhizhong Han, Mingyang Shang, Zhenbao Liu, Chi-Man Vong, Yu-Shen Liu, Matthias Zwicker, Junwei Han, and CL Philip Chen. Seqviews2seqlabels: Learning 3d global features via aggregating sequential views by rnn with attention. *IEEE Transactions on Image Processing*, 28(2):658–672, 2018.

- [40] Zhizhong Han, Honglei Lu, Zhenbao Liu, Chi-Man Vong, Yu-Shen Liu, Matthias Zwicker, Junwei Han, and CL Philip Chen. 3d2seqviews: Aggregating sequential views for 3d global feature learning by cnn with hierarchical attention aggregation. *IEEE Transactions on Image Processing*, 28(8):3986–3999, 2019.
- [41] Xin Wei, Ruixuan Yu, and Jian Sun. View-gcn: View-based graph convolutional network for 3d shape analysis. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1850–1859, 2020.
- [42] Yinan Wang, Wenbo Sun, Jionghua Jin, Zhenyu Kong, and Xiaowei Yue. Mvgcn: Multi-view graph convolutional neural network for surface defect identification using three-dimensional point cloud. *Journal of Manufacturing Science and Engineering*, 145(3):031004, 2023.
- [43] Qi Liang, Qiang Li, Lihu Zhang, Haixiao Mi, Weizhi Nie, and Xuanya Li. Mhfp: Multi-view based hierarchical fusion pooling method for 3d shape recognition. *Pattern Recognition Letters*, 150:214–220, 2021.
- [44] Wenju Wang, Xiaolin Wang, Gang Chen, and Haoran Zhou. Multi-view softpool attention convolutional networks for 3d model classification. *Frontiers in Neuroinformatics*, 16:1029968, 2022.
- [45] Daniel Maturana and Sebastian Scherer. Voxnet: A 3d convolutional neural network for real-time object recognition. In *2015 IEEE/RSJ international conference on intelligent robots and systems (IROS)*, pages 922–928. IEEE, 2015.
- [46] Yangyan Li, Soeren Pirk, Hao Su, Charles R Qi, and Leonidas J Guibas. Fpnn: Field probing neural networks for 3d data. *Advances in neural information processing systems*, 29, 2016.
- [47] Gernot Riegler, Ali Osman Ulusoy, and Andreas Geiger. Octnet: Learning deep 3d representations at high resolutions. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3577–3586, 2017.
- [48] Roman Klokov and Victor Lempitsky. Escape from cells: Deep kd-networks for the recognition of 3d point cloud models. In *Proceedings of the IEEE international conference on computer vision*, pages 863–872, 2017.
- [49] Yizhak Ben-Shabat, Michael Lindenbaum, and Anath Fischer. 3dmfv: Three-dimensional point cloud classification in real-time using convolutional neural networks. *IEEE Robotics and Automation Letters*, 3(4):3145–3152, 2018.
- [50] Truc Le and Ye Duan. Pointgrid: A deep network for 3d shape understanding. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 9204–9214, 2018.
- [51] AAM Muzahid, Wanggen Wan, Ferdous Sohel, Naimat Ullah Khan, Ofelia Delfina Cervantes Villagómez, and Hidayat Ullah. 3d object classification using a volumetric deep neural network: An efficient octree guided auxiliary learning approach. *IEEE Access*, 8:23802–23816, 2020.
- [52] Charles Ruizhongtai Qi, Li Yi, Hao Su, and Leonidas J Guibas. Pointnet++: Deep hierarchical feature learning on point sets in a metric space. *Advances in neural information processing systems*, 30, 2017.
- [53] Majid Nasrollahi, Neshat Bolourian, and Amin Hammad. Concrete surface defect detection using deep neural network based on lidar scanning. In *Proceedings of the CSCE Annual Conference, Laval, Greater Montreal, QC, Canada*, pages 12–15, 2019.
- [54] Jin-Man Park, Yong-Ho Yoo, Ue-Hwan Kim, Dukyoung Lee, and Jong-Hwan Kim. D 3 pointnet: Dual-level defect detection pointnet for solder paste printer in surface mount technology. *IEEE Access*, 8:140310–140322, 2020.
- [55] Jiaxin Li, Ben M Chen, and Gim Hee Lee. So-net: Self-organizing network for point cloud analysis. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 9397–9406, 2018.
- [56] Guocheng Qian, Yuchen Li, Houwen Peng, Jinjie Mai, Hasan Abed Al Kader Hammoud, Mohamed Elhoseiny, and Bernard Ghanem. Pointnext: Revisiting pointnet++ with improved training and scaling strategies, 2022.
- [57] Mor Joseph-Rivlin, Alon Zvirin, and Ron Kimmel. Momen (e) t: Flavor the moments in learning to classify shapes. In *Proceedings of the IEEE/CVF international conference on computer vision workshops*, pages 0–0, 2019.
- [58] Jiancheng Yang, Qiang Zhang, Bingbing Ni, Linguo Li, Jinxian Liu, Mengdie Zhou, and Qi Tian. Modeling point clouds with self-attention and gumbel subset sampling. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 3323–3332, 2019.
- [59] Hengshuang Zhao, Li Jiang, Chi-Wing Fu, and Jiaya Jia. Pointweb: Enhancing local neighborhood features for point cloud processing. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 5565–5573, 2019.

- [60] Juan Du, Hao Yan, Tzyy-Shuh Chang, and Jianjun Shi. A tensor voting-based surface anomaly classification approach by using 3d point cloud data. *Journal of Manufacturing Science and Engineering*, 144(5):051005, 2022.
- [61] Xiao Sun, Zhouhui Lian, and Jianguo Xiao. Srinet: Learning strictly rotation-invariant representations for point cloud classification and segmentation. In *Proceedings of the 27th ACM international conference on multimedia*, pages 980–988, 2019.
- [62] Xu Yan, Chaoda Zheng, Zhen Li, Sheng Wang, and Shuguang Cui. Pointasnl: Robust point clouds processing using nonlocal neural networks with adaptive sampling. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 5589–5598, 2020.
- [63] Amir Hertz, Rana Hanocka, Raja Giryes, and Daniel Cohen-Or. Pointgmm: A neural gmm network for point clouds. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12054–12063, 2020.
- [64] Niannian Wang, Duo Ma, Xueming Du, Bin Li, Danyang Di, Gaozhao Pang, and Yihang Duan. An automatic defect classification and segmentation method on three-dimensional point clouds for sewer pipes. *Tunnelling and Underground Space Technology*, 143:105480, 2024.
- [65] Meng-Hao Guo, Jun-Xiong Cai, Zheng-Ning Liu, Tai-Jiang Mu, Ralph R Martin, and Shi-Min Hu. Pct: Point cloud transformer. *Computational Visual Media*, 7:187–199, 2021.
- [66] Joakim Bruslund Haurum, Moaz MJ Allahham, Mathias S Lyng, Kasper Schøn Henriksen, Ivan A Nikolov, and Thomas B Moeslund. Sewer defect classification using synthetic point clouds. In *VISIGRAPP (5: VISAPP)*, pages 891–900, 2021.
- [67] Yunxiang Zhou, Ankang Ji, and Limao Zhang. Sewer defect detection from 3d point clouds using a transformer-based deep learning model. *Automation in Construction*, 136:104163, 2022.
- [68] Varun Kasireddy and Burcu Akinci. Encoding 3d point contexts for self-supervised spall classification using 3d bridge point clouds. *Journal of Computing in Civil Engineering*, 37(2):04022061, 2023.
- [69] Hengshuang Zhao, Li Jiang, Jiaya Jia, Philip Torr, and Vladlen Koltun. Point transformer, 2021.
- [70] Luyao Liu, Enqing Chen, and Yingqiang Ding. Tr-net: a transformer-based neural network for point cloud processing. *Machines*, 10(7):517, 2022.
- [71] Xumin Yu, Lulu Tang, Yongming Rao, Tiejun Huang, Jie Zhou, and Jiwen Lu. Point-bert: Pre-training 3d point cloud transformers with masked point modeling. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 19313–19322, 2022.
- [72] Yatian Pang, Wenxiao Wang, Francis EH Tay, Wei Liu, Yonghong Tian, and Li Yuan. Masked autoencoders for point cloud self-supervised learning. In *European conference on computer vision*, pages 604–621. Springer, 2022.
- [73] Yahui Liu, Bin Tian, Yisheng Lv, Lingxi Li, and Fei-Yue Wang. Point cloud classification using content-based transformer via clustering in feature space. *IEEE/CAA Journal of Automatica Sinica*, 2023.
- [74] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. *Advances in neural information processing systems*, 25, 2012.
- [75] Yann LeCun, Bernhard Boser, John Denker, Donnie Henderson, Richard Howard, Wayne Hubbard, and Lawrence Jackel. Handwritten digit recognition with a back-propagation network. *Advances in neural information processing systems*, 2, 1989.
- [76] Pedro Hermosilla, Tobias Ritschel, Pere-Pau Vázquez, Àlvar Vinacua, and Timo Ropinski. Monte carlo convolution for learning on non-uniformly sampled point clouds. *ACM Transactions on Graphics (TOG)*, 37(6): 1–12, 2018.
- [77] Yifan Xu, Tianqi Fan, Mingye Xu, Long Zeng, and Yu Qiao. Spidercnn: Deep learning on point sets with parameterized convolutional filters. In *Proceedings of the European conference on computer vision (ECCV)*, pages 87–102, 2018.
- [78] Alexandre Boulch, Gilles Puy, and Renaud Marlet. Fkaconv: Feature-kernel alignment for point cloud convolution. In *Proceedings of the Asian Conference on Computer Vision*, 2020.
- [79] Wenxuan Wu, Zhongang Qi, and Li Fuxin. Pointconv: Deep convolutional networks on 3d point clouds. In *Proceedings of the IEEE/CVF Conference on computer vision and pattern recognition*, pages 9621–9630, 2019.
- [80] Hugues Thomas, Charles R Qi, Jean-Emmanuel Deschaud, Beatriz Marcotegui, François Goulette, and Leonidas J Guibas. Kpconv: Flexible and deformable convolution for point clouds. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 6411–6420, 2019.

- [81] Alexandre Boulch. Convpoint: Continuous convolutions for point cloud processing. *Computers & Graphics*, 88: 24–34, 2020.
- [82] Yangyan Li, Rui Bu, Mingchao Sun, Wei Wu, Xinhan Di, and Baoquan Chen. Pointcnn: Convolution on x-transformed points. *Advances in neural information processing systems*, 31, 2018.
- [83] Binh-Son Hua, Minh-Khoi Tran, and Sai-Kit Yeung. Pointwise convolutional neural networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 984–993, 2018.
- [84] Adrien Poulenard, Marie-Julie Rakotosaona, Yann Ponty, and Maks Ovsjanikov. Effective rotation-invariant point cnn with spherical harmonics kernels. In *2019 International Conference on 3D Vision (3DV)*, pages 47–56. IEEE, 2019.
- [85] Changshuo Wang, Xin Ning, Linjun Sun, Liping Zhang, Weijun Li, and Xiao Bai. Learning discriminative features by covering local geometric space for point cloud analysis. *IEEE Transactions on Geoscience and Remote Sensing*, 60:1–15, 2022.
- [86] Martin Simonovsky and Nikos Komodakis. Dynamic edge-conditioned filters in convolutional neural networks on graphs. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3693–3702, 2017.
- [87] Yingxue Zhang and Michael Rabbat. A graph-cnn for 3d point cloud classification. In *2018 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 6279–6283. IEEE, 2018.
- [88] Jinxian Liu, Bingbing Ni, Caiyuan Li, Jiancheng Yang, and Qi Tian. Dynamic points agglomeration for hierarchical point sets learning. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 7546–7555, 2019.
- [89] Kuangen Zhang, Ming Hao, Jing Wang, Xinxing Chen, Yuquan Leng, Clarence W de Silva, and Chenglong Fu. Linked dynamic graph cnn: Learning through point cloud by linking hierarchical features. In *2021 27th international conference on mechatronics and machine vision in practice (M2VIP)*, pages 7–12. IEEE, 2021.
- [90] Seyed Saber Mohammadi, Yiming Wang, and Alessio Del Bue. Pointview-gcn: 3d shape classification with multi-view point clouds. In *2021 IEEE International Conference on Image Processing (ICIP)*, pages 3103–3107. IEEE, 2021.
- [91] Long Hoang, Suk-Hwan Lee, Eung-Joo Lee, and Ki-Ryong Kwon. Gsv-net: A multi-modal deep learning network for 3d point cloud classification. *Applied Sciences*, 12(1):483, 2022.
- [92] Lifang Chen and Qian Zhang. Ddgcnn: graph convolution network based on direction and distance for point cloud learning. *The Visual Computer*, 39(3):863–873, 2023.
- [93] Jiaying Zhang, Xiaoli Zhao, Zheng Chen, and Zhejun Lu. A review of deep learning-based semantic segmentation for point cloud. *IEEE access*, 7:179118–179133, 2019.
- [94] Alok Jhaldiyal and Navendu Chaudhary. Semantic segmentation of 3d lidar data using deep learning: a review of projection-based methods. *Applied Intelligence*, 53(6):6844–6855, 2023.
- [95] Tian Xia, Jian Yang, and Long Chen. Automated semantic segmentation of bridge point cloud based on local descriptor and machine learning. *Automation in Construction*, 133:103992, 2022.
- [96] Felix Järemo Lawin, Martin Danelljan, Patrik Tosteberg, Goutam Bhat, Fahad Shahbaz Khan, and Michael Felsberg. Deep projective 3d semantic segmentation. In *Computer Analysis of Images and Patterns: 17th International Conference, CAIP 2017, Ystad, Sweden, August 22-24, 2017, Proceedings, Part I 17*, pages 95–107. Springer, 2017.
- [97] Alexandre Boulch, Bertrand Le Saux, Nicolas Audebert, et al. Unstructured point cloud semantic labeling using deep segmentation networks. *3dor@ eurographics*, 3:17–24, 2017.
- [98] Vijay Badrinarayanan, Alex Kendall, and Roberto Cipolla. Segnet: A deep convolutional encoder-decoder architecture for image segmentation. *IEEE transactions on pattern analysis and machine intelligence*, 39(12): 2481–2495, 2017.
- [99] Nicolas Audebert, Bertrand Le Saux, and Sébastien Lefèvre. Semantic segmentation of earth observation data using multimodal and multi-scale deep networks. In *Computer Vision—ACCV 2016: 13th Asian Conference on Computer Vision, Taipei, Taiwan, November 20-24, 2016, Revised Selected Papers, Part I 13*, pages 180–196. Springer, 2017.
- [100] Alexandre Boulch, Joris Guerry, Bertrand Le Saux, and Nicolas Audebert. Snapnet: 3d point cloud semantic labeling with 2d deep segmentation networks. *Computers & Graphics*, 71:189–198, 2018.

- [101] Maxim Tatarchenko, Jaesik Park, Vladlen Koltun, and Qian-Yi Zhou. Tangent convolutions for dense prediction in 3d. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3887–3896, 2018.
- [102] Haibo Qiu, Baosheng Yu, and Dacheng Tao. Gfnet: Geometric flow network for 3d point cloud semantic segmentation. *arXiv preprint arXiv:2207.02605*, 2022.
- [103] Bichen Wu, Alvin Wan, Xiangyu Yue, and Kurt Keutzer. Squeezeseg: Convolutional neural nets with recurrent crf for real-time road-object segmentation from 3d lidar point cloud. In *2018 IEEE international conference on robotics and automation (ICRA)*, pages 1887–1893. IEEE, 2018.
- [104] Bichen Wu, Xuanyu Zhou, Sicheng Zhao, Xiangyu Yue, and Kurt Keutzer. Squeezesegv2: Improved model structure and unsupervised domain adaptation for road-object segmentation from a lidar point cloud. In *2019 international conference on robotics and automation (ICRA)*, pages 4376–4382. IEEE, 2019.
- [105] Andres Milioto, Ignacio Vizzo, Jens Behley, and Cyrill Stachniss. Rangenet++: Fast and accurate lidar semantic segmentation. In *2019 IEEE/RSJ international conference on intelligent robots and systems (IROS)*, pages 4213–4220. IEEE, 2019.
- [106] Joseph Redmon and Ali Farhadi. Yolov3: An incremental improvement. *arXiv preprint arXiv:1804.02767*, 2018.
- [107] Jinhui Zhang, Haiyan Jiang, Huizhi Shao, Qingjun Song, Xiaoshuang Wang, and Dashuai Zong. Semantic segmentation of in-vehicle point cloud with improved rangenet++ loss function. *IEEE Access*, 11:8569–8580, 2023. doi:10.1109/ACCESS.2023.3238415.
- [108] Yin Zhou and Oncel Tuzel. Voxelnet: End-to-end learning for point cloud based 3d object detection. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4490–4499, 2018.
- [109] Jing Huang and Suya You. Point cloud labeling using 3d convolutional neural network. In *2016 23rd International Conference on Pattern Recognition (ICPR)*, pages 2670–2675. IEEE, 2016.
- [110] Mehrdad S Dizaji and Devin K Harris. 3d inspectionnet: a deep 3d convolutional neural networks based approach for 3d defect detection on concrete columns. In *Nondestructive Characterization and Monitoring of Advanced Materials, Aerospace, Civil Infrastructure, and Transportation XIII*, volume 10971, pages 67–77. SPIE, 2019.
- [111] Lyne Tchammi, Christopher Choy, Iro Armeni, JunYoung Gwak, and Silvio Savarese. Segcloud: Semantic segmentation of 3d point clouds. In *2017 international conference on 3D vision (3DV)*, pages 537–547. IEEE, 2017.
- [112] Hsien-Yu Meng, Lin Gao, Yu-Kun Lai, and Dinesh Manocha. Vv-net: Voxel vae net with group convolutions for point cloud segmentation. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 8500–8508, 2019.
- [113] Eren Erdal Aksoy, Saimir Baci, and Selcuk Cavdar. Salsanet: Fast road and vehicle segmentation in lidar point clouds for autonomous driving. In *2020 IEEE intelligent vehicles symposium (IV)*, pages 926–932. IEEE, 2020.
- [114] Tiago Cortinhal, George Tzelepis, and Eren Erdal Aksoy. Salsanext: Fast, uncertainty-aware semantic segmentation of lidar point clouds. In *Advances in Visual Computing: 15th International Symposium, ISVC 2020, San Diego, CA, USA, October 5–7, 2020, Proceedings, Part II 15*, pages 207–222. Springer, 2020.
- [115] Ran Cheng, Ryan Razani, Ehsan Taghavi, Enxu Li, and Bingbing Liu. 2-s3net: Attentive feature fusion with adaptive feature selection for sparse semantic segmentation network. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 12547–12556, 2021.
- [116] Ran Cheng, Ryan Razani, Yuan Ren, and Liu Bingbing. S3net: 3d lidar sparse semantic segmentation network. In *2021 IEEE International Conference on Robotics and Automation (ICRA)*, pages 14040–14046. IEEE, 2021.
- [117] Ran Cheng, Christopher Agia, Yuan Ren, Xinhai Li, and Liu Bingbing. S3cnet: A sparse semantic scene completion network for lidar point clouds. In *Conference on Robot Learning*, pages 2148–2161. PMLR, 2021.
- [118] Xiaoyan Li, Gang Zhang, Hongyu Pan, and Zhenhua Wang. Cpgnet: Cascade point-grid fusion network for real-time lidar semantic segmentation. In *2022 International Conference on Robotics and Automation (ICRA)*, pages 11117–11123. IEEE, 2022.
- [119] Jaehyun Park, Chansoo Kim, Soyeong Kim, and Kichun Jo. Pscsnet: Fast 3d semantic segmentation of lidar point cloud for autonomous car using point convolution and sparse convolution network. *Expert Systems with Applications*, 212:118815, 2023.
- [120] Miguel Martin-Abadal, Manuel Piñar-Molina, Antoni Martorell-Torres, Gabriel Oliver-Codina, and Yolanda Gonzalez-Cid. Underwater pipe and valve 3d recognition using deep learning segmentation. *Journal of Marine Science and Engineering*, 9(1):5, 2020.

- [121] Neshat Bolourian, Majid Nasrollahi, Fardin Bahreini, and Amin Hammad. Point cloud-based concrete surface defect semantic segmentation. *Journal of Computing in Civil Engineering*, 37(2):04022056, 2023.
- [122] Jiangpeng Shu, Wenhao Li, Congguang Zhang, Yifan Gao, Yiqiang Xiang, and Ling Ma. Point cloud-based dimensional quality assessment of precast concrete components using deep learning. *Journal of Building Engineering*, 70:106391, 2023.
- [123] Chao Yin, Boyu Wang, Vincent JL Gan, Mingzhu Wang, and Jack CP Cheng. Automated semantic segmentation of industrial point clouds using respointnet++. *Automation in Construction*, 130:103874, 2021.
- [124] Mingyang Jiang, Yiran Wu, Tianqi Zhao, Zelin Zhao, and Cewu Lu. Pointsift: A sift-like network module for 3d point cloud semantic segmentation. *arXiv preprint arXiv:1807.00652*, 2018.
- [125] Francis Engelmann, Theodora Kontogianni, Jonas Schult, and Bastian Leibe. Know what your neighbors do: 3d semantic segmentation of point clouds. In *Proceedings of the European Conference on Computer Vision (ECCV) Workshops*, pages 0–0, 2018.
- [126] Qingyong Hu, Bo Yang, Linhai Xie, Stefano Rosa, Yulan Guo, Zhihua Wang, Niki Trigoni, and Andrew Markham. Randla-net: Efficient semantic segmentation of large-scale point clouds. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 11108–11117, 2020.
- [127] Chuanyu Luo, Xiaohan Li, Nuo Cheng, Han Li, Shengguang Lei, and Pu Li. Mvp-net: Multiple view pointwise semantic segmentation of large-scale point clouds. *arXiv preprint arXiv:2201.12769*, 2022.
- [128] Varun Kasireddy and Burcu Akinci. Assessing the impact of 3d point neighborhood size selection on unsupervised spall classification with 3d bridge point clouds. *Advanced Engineering Informatics*, 52:101624, 2022.
- [129] Lin-Zhuo Chen, Xuan-Yi Li, Deng-Ping Fan, Kai Wang, Shao-Ping Lu, and Ming-Ming Cheng. Lsanet: Feature learning on point sets by local spatial aware layer. *arXiv preprint arXiv:1905.05442*, 2019.
- [130] Chenxi Zhao, Weihao Zhou, Li Lu, and Qijun Zhao. Pooling scores of neighboring points for improved 3d point cloud segmentation. In *2019 IEEE international conference on image processing (ICIP)*, pages 1475–1479. IEEE, 2019.
- [131] Zhiyu Hu, Dongbo Zhang, Shuai Li, and Hong Qin. Attention-based relation and context modeling for point cloud semantic segmentation. *Computers & Graphics*, 90:126–134, 2020.
- [132] Shuang Deng and Qiulei Dong. Ga-net: Global attention network for point cloud semantic segmentation. *IEEE Signal Processing Letters*, 28:1300–1304, 2021.
- [133] Ankang Ji, Yunxiang Zhou, Limao Zhang, Robert LK Tiong, and Xiaolong Xue. Semi-supervised learning-based point cloud network for segmentation of 3d tunnel scenes. *Automation in Construction*, 146:104668, 2023.
- [134] Yunxiang Zhou, Ankang Ji, Limao Zhang, and Xiaolong Xue. Attention-enhanced sampling point cloud network (aspcnet) for efficient 3d tunnel semantic segmentation. *Automation in Construction*, 146:104667, 2023.
- [135] Jaesung Choe, Chunghyun Park, Francois Rameau, Jaesik Park, and In So Kweon. Pointmixer: Mlp-mixer for point cloud understanding. In *European Conference on Computer Vision*, pages 620–640. Springer, 2022.
- [136] Na Zhao, Tat-Seng Chua, and Gim Hee Lee. Ps²-net: A locally and globally aware network for point-based semantic segmentation. In *2020 25th International Conference on Pattern Recognition (ICPR)*, pages 723–730. IEEE, 2021.
- [137] Relja Arandjelovic, Petr Gronat, Akihiko Torii, Tomas Pajdla, and Josef Sivic. Netvlad: Cnn architecture for weakly supervised place recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 5297–5307, 2016.
- [138] Siqi Fan, Qiulei Dong, Fenghua Zhu, Yisheng Lv, Peijun Ye, and Fei-Yue Wang. Scf-net: Learning spatial contextual features for large-scale point cloud segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14504–14513, 2021.
- [139] Yuanwei Bi, Lujian Zhang, Yaowen Liu, Yansen Huang, and Hao Liu. A local-global feature fusing method for point clouds semantic segmentation. *IEEE Access*, 2023.
- [140] Yiqiang Zhao, Xingyi Ma, Bin Hu, Qi Zhang, Mao Ye, and Guoqing Zhou. A large-scale point cloud semantic segmentation network via local dual features and global correlations. *Computers & Graphics*, 111:133–144, 2023.
- [141] Xiaohan Tu, Chuanhao Zhang, Siping Liu, Cheng Xu, and Renfa Li. Point cloud segmentation of overhead contact systems with deep learning in high-speed rails. *Journal of Network and Computer Applications*, 216: 103671, 2023.

- [142] Anh-Thuan Tran, Hoanh-Su Le, Suk-Hwan Lee, and Ki-Ryong Kwon. Pointct: Point central transformer network for weakly-supervised point cloud semantic segmentation. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 3556–3565, 2024.
- [143] Junfeng Jing and Huaqing Wang. Defect segmentation with local embedding in industrial 3d point clouds based on transformer. *Measurement Science and Technology*, 35(3):035406, 2023.
- [144] Yeritza Perez-Perez, Mani Golparvar-Fard, and Khaled El-Rayes. Scan2bim-net: Deep learning method for segmentation of point clouds for scan-to-bim. *Journal of Construction Engineering and Management*, 147(9): 04021107, 2021.
- [145] Zhenxing Xu, Aizeng Wang, Fei Hou, and Gang Zhao. Defect detection of gear parts in virtual manufacturing. *Visual Computing for Industry, Biomedicine, and Art*, 6(1):1–12, 2023.
- [146] Shenlong Wang, Simon Suo, Wei-Chiu Ma, Andrei Pokrovsky, and Raquel Urtasun. Deep parametric continuous convolutional neural networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2589–2597, 2018.
- [147] Yong Li, Xu Li, Zhenxin Zhang, Feng Shuang, Qi Lin, and Jincheng Jiang. Densekpnet: Dense kernel point convolutional neural networks for point cloud semantic segmentation. *IEEE Transactions on Geoscience and Remote Sensing*, 60:1–13, 2022.
- [148] Guofeng Tong, Yuyuan Shao, and Hao Peng. Learning local contextual features for 3d point clouds semantic segmentation by attentive kernel convolution. *The Visual Computer*, pages 1–17, 2023.
- [149] Francis Engelmann, Theodora Kontogianni, and Bastian Leibe. Dilated point convolutions: On the receptive field size of point convolutions on 3d point clouds. In *2020 IEEE International Conference on Robotics and Automation (ICRA)*, pages 9463–9469. IEEE, 2020.
- [150] Wenxuan Wu, Li Fuxin, and Qi Shan. Pointconvformer: Revenge of the point-based convolution. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 21802–21813, 2023.
- [151] F Bahreini and A Hammad. Point cloud semantic segmentation of concrete surface defects using dynamic graph cnn. In *ISARC. Proceedings of the International Symposium on Automation and Robotics in Construction*, volume 38, pages 379–386. IAARC Publications, 2021.
- [152] Hyeonsoo Kim and Changwan Kim. Deep-learning-based classification of point clouds for bridge inspection. *Remote Sensing*, 12(22):3757, 2020.
- [153] Roberto Pierdicca, Marina Paolanti, Francesca Matrone, Massimo Martini, Christian Morbidoni, Eva Savina Malinverni, Emanuele Frontoni, and Andrea Maria Lingua. Point cloud semantic segmentation using a deep learning framework for cultural heritage. *Remote Sensing*, 12(6):1005, 2020.
- [154] Jie Nie, Lei Huang, Chengyu Zheng, Xiaowei Lv, and Rui Wang. Cross-scale graph interaction network for semantic segmentation of remote sensing images. *ACM Transactions on Multimedia Computing, Communications and Applications*, 19(6):1–18, 2023.
- [155] Chengzhi Wu, Xuelei Bi, Julius Pfrommer, Alexander Cebulla, Simon Mangold, and Jürgen Beyerer. Sim2real transfer learning for point cloud segmentation: An industrial application case on autonomous disassembly. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 4531–4540, 2023.
- [156] Mingtao Feng, Liang Zhang, Xuefei Lin, Syed Zulqarnain Gilani, and Ajmal Mian. Point attention network for semantic segmentation of 3d point clouds. *Pattern Recognition*, 107:107446, 2020.
- [157] Zijin Du, Hailiang Ye, and Feilong Cao. A novel local-global graph convolutional method for point cloud semantic segmentation. *IEEE Transactions on Neural Networks and Learning Systems*, 2022.
- [158] In-Ho Kim, Haemin Jeon, Seung-Chan Baek, Won-Hwa Hong, and Hyung-Jo Jung. Application of crack identification techniques for an aging concrete bridge inspection using an unmanned aerial vehicle. *Sensors*, 18(6):1881, 2018.
- [159] Javier Grandio, Belen Riveiro, Daniel Lamas, and Pedro Arias. Multimodal deep learning for point cloud panoptic segmentation of railway environments. *Automation in Construction*, 150:104854, 2023.
- [160] Zhaoxiang Zhang, Ankang Ji, Limao Zhang, Yuelei Xu, and Qing Zhou. Deep learning for large-scale point cloud segmentation in tunnels considering causal inference. *Automation in Construction*, 152:104915, 2023.
- [161] Eva Agapaki and Ioannis Brilakis. Instance segmentation of industrial point cloud data. *Journal of Computing in Civil Engineering*, 35(6):04021022, 2021.
- [162] Helin Li, Bin Lin, Chen Zhang, Liang Xu, Tianyi Sui, Yang Wang, Xinquan Hao, Deyu Lou, and Hongyu Li. Mask-point: automatic 3d surface defects detection network for fiber-reinforced resin matrix composites. *Polymers*, 14(16):3390, 2022.

- [163] Ji Hou, Angela Dai, and Matthias Nießner. 3d-sis: 3d semantic instance segmentation of rgb-d scans. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 4421–4430, 2019.
- [164] Jinhua Lin, Lin Ma, and Yu Yao. Segmentation of casting defect regions for the extraction of microstructural properties. *Engineering applications of artificial intelligence*, 85:150–163, 2019.
- [165] 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.
- [166] Ze Yang and Liwei Wang. Learning relationships for multi-view 3d object recognition. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 7505–7514, 2019.
- [167] Shih-Hung Liu, Shang-Yi Yu, Shao-Chi Wu, Hwann-Tzong Chen, and Tyng-Luh Liu. Learning gaussian instance segmentation in point clouds. *arXiv preprint arXiv:2007.09860*, 2020.
- [168] Lei Han, Tian Zheng, Lan Xu, and Lu Fang. Occuseg: Occupancy-aware 3d instance segmentation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 2940–2949, 2020.
- [169] Ze Liu, Zheng Zhang, Yue Cao, Han Hu, and Xin Tong. Group-free 3d object detection via transformers. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 2949–2958, 2021.
- [170] Ruohao Guo, Dantong Niu, Liao Qu, and Zhenbo Li. Sotr: Segmenting objects with transformers. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 7157–7166, 2021.
- [171] Yutong Xie, Jianpeng Zhang, Chunhua Shen, and Yong Xia. Cotr: Efficiently bridging cnn and transformer for 3d medical image segmentation. In *Medical Image Computing and Computer Assisted Intervention—MICCAI 2021: 24th International Conference, Strasbourg, France, September 27–October 1, 2021, Proceedings, Part III 24*, pages 171–180. Springer, 2021.
- [172] Justin Lazarow, Weijian Xu, and Zhuowen Tu. Instance segmentation with mask-supervised polygonal boundary transformers. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4382–4391, 2022.
- [173] Jiahao Sun, Chunmei Qing, Junpeng Tan, and Xiangmin Xu. Superpoint transformer for 3d scene instance segmentation. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 37(2), pages 2393–2401, 2023.
- [174] Jonas Schult, Francis Engelmann, Alexander Hermans, Or Litany, Siyu Tang, and Bastian Leibe. Mask3d: Mask transformer for 3d semantic instance segmentation. In *2023 IEEE International Conference on Robotics and Automation (ICRA)*, pages 8216–8223. IEEE, 2023.
- [175] Xin Lai, Yuhui Yuan, Ruihang Chu, Yukang Chen, Han Hu, and Jiaya Jia. Mask-attention-free transformer for 3d instance segmentation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 3693–3703, 2023.
- [176] Weiyue Wang, Ronald Yu, Qianguai 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.
- [177] Chen Liu and Yasutaka Furukawa. Masc: Multi-scale affinity with sparse convolution for 3d instance segmentation. *arXiv preprint arXiv:1902.04478*, 2019.
- [178] Zhidong Liang, Ming Yang, Hao Li, and Chunxiang Wang. 3d instance embedding learning with a structure-aware loss function for point cloud segmentation. *IEEE Robotics and Automation Letters*, 5(3):4915–4922, 2020.
- [179] Xinlong Wang, Shu Liu, Xiaoyong Shen, Chunhua Shen, and Jiaya Jia. Associatively segmenting instances and semantics in point clouds. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4096–4105, 2019.
- [180] Lin Zhao and Wenbing Tao. Jsnet: Joint instance and semantic segmentation of 3d point clouds. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34 (7), pages 12951–12958, 2020.
- [181] Lin Zhao and Wenbing Tao. Jsnet++: Dynamic filters and pointwise correlation for 3d point cloud instance and semantic segmentation. *IEEE Transactions on Circuits and Systems for Video Technology*, 33(4):1854–1867, 2022.
- [182] Francis Engelmann, Martin Bokeloh, Alireza Fathi, Bastian Leibe, and Matthias Nießner. 3d-mpa: Multi-proposal aggregation for 3d semantic instance segmentation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 9031–9040, 2020.

- [183] Zhihao Liang, Zhihao Li, Songcen Xu, Mingkui Tan, and Kui Jia. Instance segmentation in 3d scenes using semantic superpoint tree networks. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 2783–2792, 2021.
- [184] Thang Vu, Kookhoi Kim, Tung M Luu, Thanh Nguyen, and Chang D Yoo. Softgroup for 3d instance segmentation on point clouds. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2708–2717, 2022.
- [185] Jean Lahoud, Bernard Ghanem, Marc Pollefeys, and Martin R Oswald. 3d instance segmentation via multi-task metric learning. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 9256–9266, 2019.
- [186] Li Jiang, Hengshuang Zhao, Shaoshuai Shi, Shu Liu, Chi-Wing Fu, and Jiaya Jia. Pointgroup: Dual-set point grouping for 3d instance segmentation. In *Proceedings of the IEEE/CVF conference on computer vision and Pattern recognition*, pages 4867–4876, 2020.
- [187] Shaoyu Chen, Jiemin Fang, Qian Zhang, Wenyu Liu, and Xinggong Wang. Hierarchical aggregation for 3d instance segmentation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 15467–15476, 2021.
- [188] Tong He, Chunhua Shen, and Anton Van Den Hengel. Dyco3d: Robust instance segmentation of 3d point clouds through dynamic convolution. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 354–363, 2021.
- [189] Zongji Wang and Feng Lu. Voxsegnet: Volumetric cnns for semantic part segmentation of 3d shapes. *IEEE transactions on visualization and computer graphics*, 26(9):2919–2930, 2019.
- [190] 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.
- [191] Evangelos Kalogerakis, Melinos Averkiou, Subhransu Maji, and Siddhartha Chaudhuri. 3d shape segmentation with projective convolutional networks. In *proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3779–3788, 2017.
- [192] Li Yi, Hao Su, Xingwen Guo, and Leonidas J Guibas. Syncspeccnn: Synchronized spectral cnn for 3d shape segmentation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2282–2290, 2017.
- [193] Pengyu Wang, Yuan Gan, Panpan Shui, Fenggen Yu, Yan Zhang, Songle Chen, and Zhengxing Sun. 3d shape segmentation via shape fully convolutional networks. *Computers & Graphics*, 76:182–192, 2018.
- [194] Chungang Zhuang, Zhe Wang, Heng Zhao, and Han Ding. Semantic part segmentation method based 3d object pose estimation with rgb-d images for bin-picking. *Robotics and Computer-Integrated Manufacturing*, 68: 102086, 2021.
- [195] Chenyang Zhu, Kai Xu, Siddhartha Chaudhuri, Li Yi, Leonidas J Guibas, and Hao Zhang. Adacoseg: Adaptive shape co-segmentation with group consistency loss. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8543–8552, 2020.
- [196] Zhiqin Chen, Kangxue Yin, Matthew Fisher, Siddhartha Chaudhuri, and Hao Zhang. Bae-net: Branched autoencoder for shape co-segmentation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 8490–8499, 2019.
- [197] Shuailin Li, Chuyu Zhang, and Xuming He. Shape-aware semi-supervised 3d semantic segmentation for medical images. In *Medical Image Computing and Computer Assisted Intervention–MICCAI 2020: 23rd International Conference, Lima, Peru, October 4–8, 2020, Proceedings, Part I 23*, pages 552–561. Springer, 2020.
- [198] Panos Achlioptas, Judy Fan, Robert Hawkins, Noah Goodman, and Leonidas J Guibas. Shapeglot: Learning language for shape differentiation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 8938–8947, 2019.
- [199] Juil Koo, Ian Huang, Panos Achlioptas, Leonidas J Guibas, and Minhyuk Sung. Partglot: Learning shape part segmentation from language reference games. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 16505–16514, 2022.
- [200] Zhenyu Shu, Sipeng Yang, Haoyu Wu, Shiqing Xin, Chaoyi Pang, Ladislav Kavan, and Ligang Liu. 3d shape segmentation using soft density peak clustering and semi-supervised learning. *Computer-Aided Design*, 145: 103181, 2022.
- [201] Ahmed Abdelreheem, Ivan Skorokhodov, Maks Ovsjanikov, and Peter Wonka. Satr: Zero-shot semantic segmentation of 3d shapes. *arXiv preprint arXiv:2304.04909*, 2023.