



IJIRCCCE

e-ISSN: 2320-9801 | p-ISSN: 2320-9798



INTERNATIONAL JOURNAL OF INNOVATIVE RESEARCH

IN COMPUTER & COMMUNICATION ENGINEERING

Volume 10, Issue 5, May 2022

ISSN INTERNATIONAL
STANDARD
SERIAL
NUMBER
INDIA

Impact Factor: 8.165



9940 572 462



6381 907 438



ijircce@gmail.com



www.ijircce.com

Advanced End-To-End Transformer Model for 3D-Object Detection

Dr.G.Arun Sampaul Thomas, B.Manasa, K.Bhanu, P.Akhila

Associate Professor, Dept. of Computer Science and Engineering, JB Institute of Engineering and Technology,
Telangana, India

Student, Dept. of Computer Science and Engineering, JB Institute of Engineering and Technology, Telangana, India

Student, Dept. of Computer Science and Engineering, JB Institute of Engineering and Technology, Telangana, India

Student, Dept. of Computer Science and Engineering, JB Institute of Engineering and Technology, Telangana, India

ABSTRACT: 3DETR is an end-to-end Transformer-based object identification approach for 3D point clouds that we propose. In comparison to previous detection approaches that use a variety of 3D-specific inductive biases, 3DETR requires just minor changes to the vanilla Transformer block. We discover that a typical Transformer with non-parametric queries and Fourier positional embedding's competes with customised topologies that use libraries of 3D-specific operators with hand-tuned hyper parameters. Nonetheless, 3DETR is theoretically simple and straightforward to implement, allowing for additional advancements by adding 3D domain knowledge. Extensive trials reveal that 3DETR surpasses the well-established and well-tuned VoteNet baselines by 9.5 percent on the difficult ScanNetV2 dataset. Furthermore, we demonstrate that 3DETR is adaptable to 3D tasks other than detection and may be used as a foundation for future study.

KEYWORDS: Transformer model, 3D-Object detection, Object detection, End-to-End Object detection, 3D-Image detecting transformer model.

I. INTRODUCTION

The purpose of 3-D item detection is to understand and discover objects in 3-D situations. These images, which can be regularly represented as factor clouds, function an unordered, sparse, and abnormal set of factors accumulated with an intensity scanner. Because in their set like character, factor clouds vary substantially from well-known grid-like visible statistics along with images and movies. Other 3-D representations, along with a couple of views [60], voxels [1], or meshes [8], require greater post-processing to create and regularly lose data thanks to quantization. As a result, factor clouds have turn out to be a not unusual place 3-D representation, spurring the introduction of specialised 3-D structures. Many cutting-edge 3-D detection techniques construct bounding containers immediately from 3-D factors. VoteNet [42], in particular, depicts 3-D detection as a set-to-set issue, that is, converting an unordered series of inputs (factor cloud) into an unordered set of outputs (bounding containers). VoteNet employs an encoder-decoder design, with the encoder being a PointNet++ network [44] that turns the unordered Input Point Cloud Decoder Attention Detections factor set into a group of factor features. The factor traits are then despatched right into a decoder, which generates the 3-D bounding containers. While successful, such designs took years to assemble via way of means of hand-encoding inductive biases, radii, and growing precise 3-D operators and loss functions. Set-to-set encoder-decoder models have developed as a competitive technique to represent 2D object identification in tandem with 3D. The latest DETR [4] model, which is based on Transformer [68], depicts 2D object identification as a set-to-set issue. Transformers' self-attention operation is meant to be permutation-invariant and capture long-range contexts, making them an excellent option for processing unordered 3D point cloud data. Inspired by this discovery, we ask: can we use Transformers to develop a 3D object detector without depending on hand-crafted inductive biases? 3DETR eliminates several hard-coded design decisions in VoteNet and PointNet++ while being straightforward to implement and comprehend. 3DETR, unlike DETR, does not use a ConvNet backbone and instead depends only on Transformers taught from scratch. Our transformer-based detection pipeline is adaptable, and any component, much like in VoteNet, may be substituted by other existing modules. Finally, we illustrate how 3D specific inductive biases can be simply integrated into 3DETR to boost its performance even more. We obtain 65.0 percent AP and 59.0 percent AP on two typical indoor 3D detection benchmarks, ScanNetV2 and SUN RGB-D, respectively, exceeding an enhanced VoteNet baseline by 9.5 percent AP50 on ScanNetV2.

II. RELATED WORK

We present a Transformer-based 3D object detection model. We expand on previous work on 3D architectures, detection, and Transformers.

Grid-based 3D Architectures: Convolution networks can be used on irregular 3D data that has been converted into regular grids. Projection techniques [3, 19, 25, 26, 59, 60, and 65] turn 3D data into 2D planes and grids. Voxelization may also be used to turn 3D data into a volumetric 3D grid [1, 12, 15, 28, 35, 49, 56, and 66]. We employ 3D point clouds directly since they are appropriate for set-based structures like the transformer.

Point cloud Architectures: Unordered point clouds are frequently acquired by 3D sensors. It is preferable to acquire permutation invariant features while utilising unordered point clouds as input. Point-wise MLP-based architectures [17, 83] such as PointNet [44] and PointNet++ [45] learn effective representations using permutation equivariant set aggregation (down sampling) and point wise MLPs. To make the number of input points in our model manageable, we apply a single down sampling procedure from [45]. Graph-based models [27, 73] can work with unsorted 3D data. DGCNN [77] and PointWeb [90] employ local neighbourhoods of points to form graphs, SPG [24] uses attribute and context similarity, while Jiang et al. [18] use point-edge interactions. Finally, architectures based on continuous point convolution may operate on point clouds. Polynomial functions, as in SpiderCNN [80], or linear functions, as in Flex Convolutions [13], can be used to construct continuous weights. Soft-assignment matrices [69] and specified ordering [28] can also be used to apply convolutions. PointConv [78] and KPConv [67] produce convolutional weights dynamically depending on the input point coordinates, whereas InterpCNN [34] interpolates weights using these coordinates. We build on the Transformer [68], which is useful for sets but not designed for 3D.

3D Object Detection: How to predict a 3D bounding box from 3D input data has been extensively investigated [23, 41, 43, 52, 54, 55, 70, 72, 93]. Many approaches use 2D projection to avoid costly 3D processes. MV3D [6] and VoxelNet [92] are a combination of 3D and 2D convolutions. Yang et al. [81] Simplifies the 3D process, [82] uses 2D projection, and [76] uses voxel "columns". Focuses on an approach that uses 3D point clouds directly [40, 51, 75, 85]. PointRCNN [51] and PVRCNN [50] are two-step recognition pipelines for 2D images that rival the popular RCNN framework [47]. These approaches are relevant to our research, but for simplicity we will develop a one-step recognition model outlined in [11, 14, 42, 84]. VoteNet [42] uses feature sampling, grouping, and voting methods developed for 3D data to identify boxes and use Huff voting for sparse point cloud input. Voting Net serves as the basis for many follow-up projects. 3DMPA [11] combines polling and ConvNet graphs to improve article suggestions and aggregated perception using specially created 3D geometric properties. HGNet [5] improves Huff voting by using a hierarchical graph network with a functional pyramid. H3DNet [89] improves VoteNet by predicting 3D primitives and using geometric loss functions. Here are some basic cognitive approaches that can be used as the basis for these advances in 3D cognition.

III. APPROACH

We briefly review prior work in 3D detection and their conceptual similarities to 3DETR.

Preliminaries: The modern day VoteNet [42] framework, that is a set-to-set prediction framework like our technique, serves as the muse for several detection fashions in 3-d. For detection, VoteNet employs a personalized 3-d encoder and decoder structure. It combines those fashions with a sparse factor cloud Hough Voting loss. The encoder is a PointNet++ [45] version that employs a combination of down sampling (set-aggregation) and up sampling (feature-propagation) operations constructed mainly for 3-d factor clouds. The VoteNet "decoder" predicts bounding containers in 3 steps: 1) every factor "votes" for a box's centre coordinate; 2) votes are aggregated inside a given radius to generate "centres"; and 3) bounding containers are expected on "centres.". BoxNet, on the opposite hand, plays extensively worse than VoteNet due to the fact balloting keeps greater statistics in sparse factor clouds and generates better 'centre' points. Theseveral hand-encoded radii utilised with inside the encoder, decoder, and loss characteristic are crucial for detection performance, as referred to with the aid of using the authors [42], and were cautiously calibrated [44, 45]. The Transformer [68] is a general-motive structure that may address constant inputs and seize massive contexts with the aid of using computing self-interest among all pairs of enter points. Both of those traits make it an high-quality candidate version for 3-d factor clouds. Following that, we display our 3DETR version, which employs a Transformer for each the encoder and decoder with little modifications and incorporates minimum hand-coded 3-d statistics. 3DETR employs a greater truthful schooling and inference process. We additionally speak the similarities and variations among the DETR version and the DETR version for 2D detection.

3DETR: Encoder-decoder Transformer: 3DETR takes a 3D point cloud as input and predicts the locations of objects as 3D bounding boxes. A point cloud is an unordered set of N points, each with its own 3-dimensional XYZ coordinates. We utilise the set aggregation down sampling procedure from [45] to down sample the points and project them to N_0 dimensional features because the number of points is relatively huge. The resultant subset of N_0 features is processed through an encoder to provide a set of N_0 features as well. Using a parallel decoding approach inspired by [4,] a decoder accepts these characteristics as input and predicts numerous bounding boxes. Both the encoder and decoder employ typical Transformer blocks marked 'pre-norm.'

Recent attempts, on the other hand, have focused on appearance-based approaches such as enhanced feature descriptors and pattern recognition algorithms. The basic idea behind these algorithms is to calculate eigenvectors from a group of vectors, each of which represents a single face picture as a raster scan vector of gray-scale pixel values. Each eigenvector, also known as an eigenface, captures a certain variance among all the vectors, and a small selection of eigenvectors captures almost all of the appearance variation of face pictures in the training set. Given a test picture represented as a vector of grayscale pixel values, its identity is established by locating the vector's nearest neighbour after it has been projected into a subspace covered by a collection of eigenvectors. Appearance-based approaches are often divided into two stages [10, 11, 12]. In the first step, a model is built from a collection of data.pictures for reference. The look of the object is included in the set. In various angles, illuminants, and lighting conditions. For example, many instances of a class of objects might exist. faces. The photos are strongly linked and can be processed quickly. compressed using, for example, the Karhunen-Loeve transformation (see PCA stands for Principal Component Analysis) [13]. In the Parts of the input picture (subimages) are used in the second phase. size as the training pictures) are optionally removed segmentation (by texture, colour, or motion) or comprehensive segmentation enumeration of image windows over the whole picture

Appearance-Based Methods: Some applications, such as bin picking, need a position estimate in 3D space (rather than only the 2D picture plane), such as when a robot must grasp specific things from an unordered bunch of objects. Typically, such applications make use of sensor systems that enable the creation of 3D data and 3D matching. Another method for determining an object's 3D pose is to estimate the projection of the object's position in 3D space onto a 2D camera picture. Many of the approaches make use of what are known as range images or depth maps, in which information about the z-direction (e.g., z-distance to the sensor) is recorded as a function of the [x,y]-position in the picture plane. This type of data representation is not "complete" 3D.

Descriptor-based Methods: When performing object identification in "real-world" settings, characterisation with geometric primitives like as lines or circular arcs is ineffective. Another consideration is that the algorithm must correct for severe backdrop clutter and occlusion, which is difficult for global appearance algorithms. Local image information assessment is necessary to deal with partial occlusion. Furthermore, when dealing with a large number of identical items or objects with smooth brightness changes, gradient-based shape information may be insufficient. To that purpose, Schmid and Mohr [16] proposed a two-stage technique for describing visual content: the first phase consists of detecting so-called interest/key spots, which are places that display some form of conspicuous property, such as a corner.

Implementation Details: We utilise PyTorch [39] to implement 3DETR and the standard nn.MultiHeadAttention module to implement the Transformer. We subsample $N_0 = 2048$ points with a single set aggregation method [45] to generate 256 dimensional point features. The 3DETR encoder is composed of three layers, each of which employs multiheaded attention with four heads and a two-layer MLP with a 'bottleneck' of 128 hidden dimensions. The 3DETR decoder comprises 8 layers and is quite similar to the encoder, with the exception that the MLP hidden dimensions are 256. In the decoder, we employ Fourier positional encodings [64] of the XYZ coordinates. MLPs for bounding box prediction have two layers and a hidden dimension of 256.

The SIFT Algorithm: SIFT is a feature identification and description technology created by Lowe [28]. (Scale Invariant Feature Transformation). This means that a picture is searched for key points. These key points

are then retrieved and represented as a vector. The generated vectors may be used to identify consistent matches between pictures for object identification, camera calibration, 3D reconstruction, and a variety of other applications [29].

SIFT is divided into three fundamental phases. First, the image's keypoints are extracted. These keypoints are then described as 128 vectors.

Adapting number of queries:As the number of queries increases, 3DETR predicts more bounding boxes, resulting in higher performance at the expense of longer running time. Our non-parametric queries in 3DETR, on the other hand, allow us to adjust the amount of box predictions to exchange performance for running time. It should be noted that while this is feasible with VoteNet, it is not with DETR. In Fig 5 (right), we compare various models trained with variable numbers of questions to modifying the number of inquiries at test time. At test time, the same 3DETR model may adapt to a varied amount of questions and performs comparable to other models. Performance improves until the number of queries is sufficient to adequately cover the point cloud.

IV. SIMULATION RESULTS

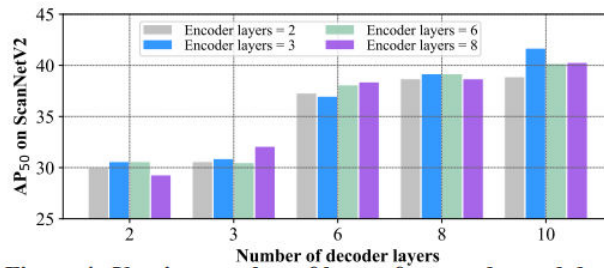
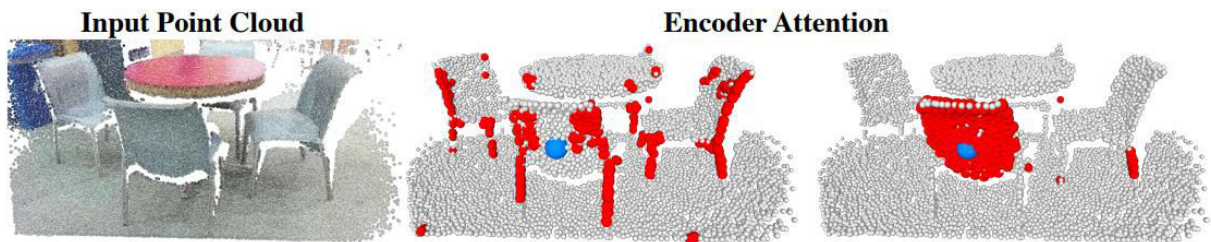
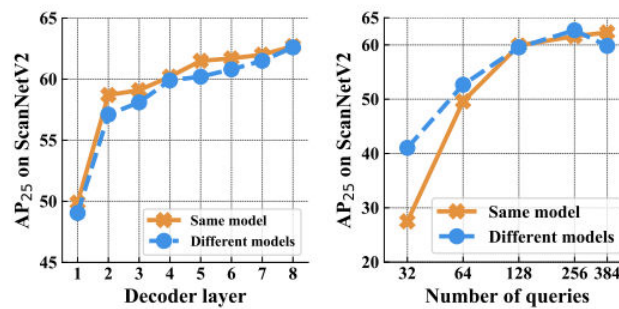


Figure 4: Varying number of layers for encoder and decoder. We train different models with varying number of encoder and decoder layers and analyze the impact on detection performance on ScanNetV2. Increasing the number of layers in either the encoder or decoder has a positive effect, but a higher number of decoder layers matters more than the encoder layers.



V. CONCLUSION AND FUTURE WORK

3DETR, an end-to-end Transformer model for 3D detection on point clouds, was presented. 3DETR needs just a few 3D-specific design considerations or hyper parameters. We demonstrate that non-parametric searches and Fourier encodings are essential for effective 3D detection performance. Our suggested design considerations allow for strong Transformers for 3D detection, as well as other 3D tasks like as shape categorization. Furthermore, our set loss function is generalizable to previous 3D designs. In summary, 3DETR is a flexible framework that can easily combine past 3D detection components and be exploited to develop more powerful 3D detectors. Finally, it incorporates the flexibility of both VoteNet and DETR, allowing for a variable number of predictions at test time (like VoteNet does) as well as a variable number of decoder layers.

REFERENCES

- [1] Andrew Adams, Jongmin Baek, and Myers Abraham Davis. Fast high-dimensional filtering using the permutohedral lattice. In Computer Graphics Forum, volume 29, pages 753–762. Wiley Online Library, 2010.
- [2] Jimmy Lei Ba, Jamie Ryan Kiros, and Geoffrey E Hinton. Layer normalization. arXiv preprint arXiv:1607.06450, 2016.
- [3] Alexandre Boulch, Bertrand Le Saux, and Nicolas Audebert. Unstructured point cloud semantic labeling using deep segmentation networks. 3DOR, 2:7, 2017.
- [4] Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and Sergey Zagoruyko. End-to-end object detection with transformers. In European Conference on Computer Vision, pages 213–229. Springer, 2020.
- [5] Jintai Chen, Biwen Lei, Qingyu Song, Haochao Ying, Danny Z Chen, and Jian Wu. A hierarchical graph network for 3d object detection on point clouds. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pages 392–401, 2020.
- [6] Xiaozhi Chen, Huimin Ma, Ji Wan, Bo Li, and Tian Xia. Multi-view 3d object detection network for autonomous driving. In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pages 1907–1915, 2017.
- [7] Angela Dai, Angel X. Chang, Manolis Savva, Maciej Halber, Thomas Funkhouser, and Matthias Niessner. Scannet: Richly-annotated 3d reconstructions of indoor scenes. In The IEEE Conference on Computer Vision and Pattern Recognition (CVPR), July 2017.
- [8] Boris Delaunay et al. Sur la sphere vide. Izv. Akad. Nauk SSSR, Otdelenie Matematicheskii i Estestvennyka Nauk, 7(793-800):1–2, 1934.
- [9] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. arXiv preprint arXiv:1810.04805, 2018.
- [10] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. arXiv preprint arXiv:2010.11929, 2020.
- [11] 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.
- [12] Ben Graham. Sparse 3d convolutional neural networks. arXiv preprint arXiv:1505.02890, 2015.
- [13] Fabian Groh, Patrick Wieschollek, and Hendrik PA Lensch. Flex-convolution. In Asian Conference on Computer Vision, pages 105–122. Springer, 2018.
- [14] JunYoung Gwak, Christopher B Choy, and Silvio Savarese. Generative sparse detection networks for 3d single-shot object detection. In European conference on computer vision, 2020.
- [15] Pedro Hermosilla, Tobias Ritschel, Pere-Pau Vazquez, Alvar 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.
- [16] Han Hu, Jiayuan Gu, Zheng Zhang, Jifeng Dai, and Yichen Wei. Relation networks for object detection. In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pages 3588–3597, 2018.
- [17] 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.
- [18] Li Jiang, Hengshuang Zhao, Shu Liu, Xiaoyong Shen, ChiWing Fu, and Jiaya Jia. Hierarchical point-edge interaction network for point cloud semantic segmentation. In Proceedings of the IEEE/CVF International Conference on Computer Vision, pages 10433–10441, 2019.
- [19] Asako Kanezaki, Yasuyuki Matsushita, and Yoshifumi Nishida. Rotationnet: Joint object categorization and pose estimation using multiviews from unsupervised viewpoints. In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pages 5010–5019, 2018.
- [20] Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In 3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings, 2015.
- [21] Guillaume Klein, Yoon Kim, Yuntian Deng, Jean Senellart, and Alexander M Rush. Opennmt: Open-source toolkit for neural machine translation. arXiv preprint arXiv:1701.02810, 2017.
- [22] Harold W Kuhn. The hungarian method for the assignment problem. Naval research logistics quarterly, 2(1-2):83–97, 1955.
- [23] Jean Lahoud, Bernard Ghanem, Marc Pollefeys, and Martin R Oswald. 3d instance segmentation via multi-task metric learning. In Proceedings of the IEEE International Conference on Computer Vision, pages 9256–9266, 2019.

- [24] Loic Landrieu and Martin Simonovsky. Large-scale point cloud semantic segmentation with superpoint graphs. In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pages 4558–4567, 2018.
- [25] Alex H Lang, Sourabh Vora, Holger Caesar, Lubing Zhou, Jiong Yang, and Oscar Beijbom. Pointpillars: Fast encoders for object detection from point clouds. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pages 12697–12705, 2019.
- [26] Felix Jaremo Lawin, Martin Danelljan, Patrik Tosteberg, Goutam Bhat, Fahad Shahbaz Khan, and Michael Felsberg. Deep projective 3d semantic segmentation. In International Conference on Computer Analysis of Images and Patterns, pages 95–107. Springer, 2017.
- [27] Guohao Li, Matthias Muller, Ali Thabet, and Bernard Ghanem. Deepgcns: Can gcns go as deep as cnns? In Proceedings of the IEEE/CVF International Conference on Computer Vision, pages 9267–9276, 2019.
- [28] Yangyan Li, Rui Bu, Mingchao Sun, Wei Wu, Xinhan Di, and Baoquan Chen. Pointcnn: Convolution on x-transformed points. In Advances in Neural Information Processing Systems, pages 820–830, 2018.
- [29] Zhe Liu, Xin Zhao, Tengpeng Huang, Ruolan Hu, Yu Zhou, and Xiang Bai. Tanet: Robust 3d object detection from point clouds with triple attention. In Proceedings of the AAAI Conference on Artificial Intelligence, volume 34, pages 11677–11684, 2020.
- [30] Ilya Loshchilov and Frank Hutter. Sgdr: Stochastic gradient descent with warm restarts. arXiv preprint arXiv:1608.03983, 2016.
- [31] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. arXiv preprint arXiv:1711.05101, 2017.
- [32] Jiasen Lu, Dhruv Batra, Devi Parikh, and Stefan Lee. Vilbert: Pretraining task-agnostic visiolinguistic representations for vision-and-language tasks. arXiv preprint arXiv:1908.02265, 2019.
- [33] Christoph Luscher, Eugen Beck, Kazuki Irie, Markus Kitza, Wilfried Michel, Albert Zeyer, Ralf Schluter, and Hermann Ney. Rwth asr systems for librispeech: Hybrid vs attention– w/o data augmentation. arXiv preprint arXiv:1905.03072, 2019.
- [34] Jiageng Mao, Xiaogang Wang, and Hongsheng Li. Interpolated convolutional networks for 3d point cloud understanding. In Proceedings of the IEEE/CVF International Conference on Computer Vision, pages 1578–1587, 2019.
- [35] 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.
- [36] Anshul Paigwar, Ozgur Erkent, Christian Wolf, and Christian Laugier. Attentional pointnet for 3d-object detection in point clouds. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops, pages 0–0, 2019.
- [37] Xuran Pan, Zhuofan Xia, Shiji Song, Li Erran Li, and Gao Huang. 3d object detection with pointformer. arXiv preprint arXiv:2012.11409, 2020.
- [38] Niki Parmar, Ashish Vaswani, Jakob Uszkoreit, Lukasz Kaiser, Noam Shazeer, Alexander Ku, and Dustin Tran. Image transformer. In International Conference on Machine Learning, pages 4055–4064. PMLR, 2018.
- [39] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Kopf, Edward Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. Pytorch: An imperative style, high-performance deep learning library. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d Alche-Buc, E. Fox, and R. Garnett, editors, Advances in Neural Information Processing Systems 32, pages 8024–8035. Curran Associates, Inc., 2019.
- [40] Quang-Hieu Pham, Thanh Nguyen, Binh-Son Hua, Gemma Roig, and Sai-Kit Yeung. Jsis3d: Joint semantic-instance segmentation of 3d point clouds with multi-task pointwise networks and multi-value conditional random fields. In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pages 8827–8836, 2019.
- [41] Trung T Pham, Markus Eich, Ian Reid, and Gordon Wyeth. Geometrically consistent plane extraction for dense indoor 3d maps segmentation. In 2016 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS), pages 4199–4204. IEEE, 2016.
- [42] Kim, Eunyoung and Gerard Medioni. "3D Object Recognition in Range Images Using Visibility Context." Intelligent Robots and Systems (IROS), 2011 IEEE/RSJ International Conference on. IEEE. (2011): 8.
- [43] Kounalakis, T. and G. A. Triantafyllidis. "3D scene's object detection and recognition using depth layers and SIFT-based machine learning." 3D Research 2.3 (2011): 11.
- [44] Bo, Liefeng, Xiaofeng Ren and Dieter Fox. "Kernel Descriptors for Visual Recognition." NIPS. Vol. 1. No. 2. (2010).
- [45] Liefeng Bo, Kevin Lai, Xiaofeng Ren Dieter Fox. "Object Recognition with Hierarchical Kernel Descriptors." Computer Vision and Pattern Recognition (CVPR), 2011 IEEE Conference on IEEE (2011): 8.



- [46] A. Torralba, R. Fergus, and W. Freeman. 80 million tiny images: A large data set for nonparametric object and scene recognition. *IEEE PAMI*, 30(11):1958–1970, 2008.
- [47] A. Krizhevsky. "Learning multiple layers of features from tiny images". Technical report, 2009.
- [48] Bo, Liefeng, Xiaofeng Ren and Dieter Fox. "Depth kernel descriptors for object recognition. "Depth Kernel Descriptors for Object Recognition." *Intelligent Robots and Systems (IROS)*, 2011 IEEE/RSJ International Conference on. IEEE (2011): 6.
- [49] Kuk-Jin Yoon, Min-Gil Shin and Ji-Hyo Lee. "Recognizing 3D Objects with 3D Information from Stereo Vision." *Pattern Recognition (ICPR)*, 2010 20th International Conference on. IEEE (2010).
- [50] Jean Ponce, Svetlana Lazebnik, Fredrick Rothganger Cordelia Schmid. "Toward True 3D Object Recognition." *Reconnaissance de Formes et Intelligence Artificielle*. (2004).
- [51] Bradski, Gary and Adrian Kaehler. *Learning OpenCV*. Ed. Mike Loukides. O'Reilly Media, 2008.
- [52] Grauman, Kristen and Bastian Leibe. "Visual Object Recognition" University of Texas at Austin and RWTH Aachen University, 2010



INNO  SPACE
SJIF Scientific Journal Impact Factor

Impact Factor: 8.165

 **doi**[®]
cross **ref**

ISSN INTERNATIONAL
STANDARD
SERIAL
NUMBER
INDIA



INTERNATIONAL JOURNAL OF INNOVATIVE RESEARCH

IN COMPUTER & COMMUNICATION ENGINEERING

 9940 572 462  6381 907 438  ijircce@gmail.com



www.ijircce.com

Scan to save the contact details