

November 25, 2020

Related work

1 Object detection overview

The crux of object detection is to locate as well as classify objects. It utilises rectangular bounding boxes in order to locate the objects that are detected and then classify the category of object. Object detection is one of the important areas of computer vision. It has several applications in scientific as well as practical industrial productions like face detection, text detection, pedestrian detection, video detection, logo detection, vehicle detection, medical image detection and so forth. The current state of the art object detection models use deep convolutional neural networks(CNNs) as their backbone as well as detection network for extracting features from input images or videos and solve the task of classification as well as localization respectively. [15]

The task of object detection can be divided into two categories:

- One stage approach
- Two-stage approach

2 One stage approach

In one stage approach class probabilities as well as bounding boxes of objects are predicted utilising a single-stage network. They don't require region proposal generation as well as post processing. So, single stage approaches are fast. [11]

The one stage approach could be divided into two categories :

- 2D object detection

- 3D object detection

2.1 2D object detection

Unified one stage 2D object detection approach refers to architectures which directly predict class probabilities as well as bounding box offsets from images with single feed-forward Convolutional Neural Network(CNN) in a monolithic setting which does not involve generation of proposal region or post classification that encapsulates all computation using a single network.

YOLO [35] divides the input image into $M \times M$ grid cells and utilizes CNNs to get the bounding box regression, confidence scores as well as class probabilities of each grid cell. YOLO0000 [36] and YOLOv3 [21] further improve the performance. Even though YOLO is fast, it misses small objects because of the coarse segmentation of input images. These drawbacks were addressed by SSD [26] by utilizing feature pyramids for single stage object detection. In SSD for every feature map locations anchor boxes of various aspect ratios and scales are generated. In RetinaNet [24] they proposed focal loss in order to handle the imbalance between target and background object bounding boxes.

2.2 3D object detection

The single-stage 3D object detection approaches [48, 54, 50, 20, 42, 22] parse the given sparse 3D point cloud to a compact representation like voxel grid or bird eye view image and use CNN to predict the bounding box. This enables single-stage approaches to become simple as well as efficient. A significant drawback of this approach is that these approaches downscale the feature maps progressively. Hence, the spatial resolution of the feature maps get lost and thus, the structural information of point cloud could not be considered explicitly. Therefore single-stage approaches are less accurate when it comes to processing the sparse point clouds.

In [21, 49], they slice the dense 3D data to 2D slices and are fed to 2D object detection models to get the prediction. In [17], Khosravan et al. proposed S4ND, a deep learning method which does lung nodule detection in one step. The main architecture is based on convolution blocks that has dense connections. They also use down-sampling methods in the network as it plays an important role in tiny object detection. In [19] they used single stage network which is implemented as 3D CNN, but it did not perform well

as it was not able to converge. Hence they used multistage detector instead. Unfortunately they didn't provide much details about the architecture of the one stage detector used.

3 Two-stage approach

In two stage approach several possible regions containing objects are proposed and then region-wise features to predict the category of each region or proposal are extracted. [11]

The two-stage approach could be divided into two categories :

- 2D object detection
- 3D object detection

3.1 2D object detection

Two-stage 2D object detection approaches are region-based frameworks. In the case of two-stage approach region proposals which are category independent are generated from an image. CNN features are then extracted from these regions. After that category specific classifiers are utilized to determine the label of the categories for the proposals.

The two-stage 2D object detection algorithms are best represented by the R-CNN family [10, 38, 13]. Faster R-CNN introduced the Region Proposal Network (RPN). A substantial number of background candidates are filtered out by RPN, and a different network is used to predict bounding box co-ordinates and class labels for each proposal. In R-FCN [6] position-sensitive feature maps are extracted. These feature maps are fed to RPN to get class scores. Mask R-CNN [14] extends Faster R-CNN to instance segmentation, they first find the bounding box coordinates and crop and segment the bounding box region to get the refined mask.

3.2 3D object detection

Two-stage 3D object detection approach [33, 39, 4, 51, 32, 40] leverage spatial information in the second stage, that focuses on the region of interest (ROI's) which are predicted by the first stage and then predicts bounding box. Compared to one stage approach two stage approach are better as

they leverage the regions of interest in first stage and focus only on those regions in second stage. This shows that accurate localization can be achieved when fine-grained spatial information is leveraged. Operating on each point and re-extracting features on each ROI increases the computational cost substantially. Hence it becomes hard for two-stage approaches to reach real-time speed. [12]

In two-stage approaches [55, 18], the dense volume data is fed to the 3D version of R-CNN models to get the prediction. In [45] slice of dense volume data, i.e. slice of 3D CT scans are fed to Mask R-CNN to get the prediction. In [52] Zhang et al. introduce pancreatic tumor detection framework which aims to fully exploit the context information at different scales. The network uses Feature Pyramid network [23] along with Faster R-CNN [37] in the backbone.

4 Datasets

Datasets have played a very important role in the history of object detection research. Datasets are not only common ground to measure and compare the performance of algorithms, but also pushed the field towards improving the system so that complex and challenging problems of object detection could be solved. In recent times deep learning techniques have been very successful in several visual recognition problems, and it could not have been possible without the availability of large amounts of annotated data.[25] The datasets for 3D object detection can be broadly classified into two categories :

- Sparse 3D point cloud
- Dense volume data

4.1 Sparse 3D point cloud

The sparse 3D point clouds could be divided into two categories namely RGB-D images and LiDAR data. In comparison to RGB-D images , LiDAR data is special. On one hand LiDAR data provides structural and spatial information of relative location and precised depth. On the other hand LiDAR data is sparse, unordered and locality sensitive, and hence it becomes more difficult to process raw LiDAR data.[51] Both data types are discussed in following sections.

4.1.1 LiDAR data

LiDAR (Light Detection and Ranging) is a remote sensing technique in which the distance to the target is gauged by illuminating the target utilizing a laser light and using a sensor in order to measure the reflection of light [34]. The difference in wavelengths and laser return times are then used to obtain a 3D representation of the target. LiDAR data has several applications in surveying, geomatics, geography, archaeology, seismology, forestry. The KITTI Vision Benchmark Suite [9] is one the datasets for 3d object detection in the field of autonomous driving. It contains LIDAR data taken using a sensor mounted in the front of the car.

4.1.2 RGB-D data

RGB-D data is a combination of RGB data along with its corresponding depth data. One of the RGB-D object detection datasets is SUN RGB-D [43].

4.2 Dense volume data

Dense volume data is also referred to as dense 3D data in [7]. Computed Tomography (CT), Magnetic Resonance Imaging (MRI) are some of the examples of dense volume data.

4.2.1 Computed Tomography (CT)

The term “Computed Tomography” refers to computerized X-ray imaging process in which narrow beam of X-rays are aimed at the patient and rotated quickly around the body which produces signals that are utilized by machine’s computer to get cross-sectional images or slices of the body. These slices are referred to as tomographic images and contains more detailed information in comparison to conventional x-rays. [29] CT scans are not limited to medical domain; they are also used in industries. In industries, they are used for detection of flaws like cracks and voids as well as particle analysis of materials. They are used in metrology for the measurement of internal and external geometry of complex parts. Some of the CT scan object detection datasets available are RSNA pneumonia detection challenge [31], COVID 19 CT scans [16], DeepLesion dataset [2], covid-19-chest-xray-lung-bounding-boxes-dataset [5]

4.2.2 Magnetic Resonance Imaging (MRI)

Magnetic Resonance Imaging(MRI) is a non-invasive imaging technique that creates three-dimensional anatomical images. It is used for diagnosis and detection of diseases as well as treatment monitoring. It relies on sophisticated technology which excites as well as detects the change in the direction of the rotational axis of protons that are found in water which makes up living tissues. MRI scanners are well suited to get the image of the non-bony parts as well as soft tissues of the body. They are different from CT scans in the way that they don’t use damaging ionizing radiation of the x-rays. The spinal cord, brain, nerves, ligaments, muscles, tendons are visible more clearly with

MRI in comparison to regular x-rays and CT, this is the reason why MRI is used for imaging the knee as well as shoulder injuries. [30]

In [53, 3], they use MRI data as an input for the task of object detection.

5 Resolution of image

Resolution means the number of pixels present in an image. Resolution is identified by the width and height of the image. Resolution also refers to the total pixels an image contains, for example, if an image is 2048 pixels wide and 1536 pixels high. Total pixels = $2048 \times 1536 = 3145728$ pixels or 3.1 Megapixels. So the resolution image is 3.1 Megapixels.[28]. One of the factors that impact the performance of object detection models is the resolution of images [27].

5.1 Advantage of using low resolution image

The advantages of using low resolution images are:

- Low memory requirement
- Frame Per Second(Fps) increases

The advantages are explained in the following sections.

5.1.1 Low memory requirement

The memory required to store the images reduces when we use the image of lower resolution; this is well illustrated in figure 1.

Inch size (changed)	Resolution (changed)	Pixel dimensions (you set)	File size
2 x 2 in	100 ppi	200 x 200 px	117.2 KB
3 x 3 in	100 ppi	300 x 300 px	263.7 KB
6 x 6 in	100 ppi	600 x 600 px	1.03 MB

Figure 1: memory requirement at various resolutions [1]

5.1.2 Frame Per Second(Fps) increases

As the resolution of the image decreases, the fps of the object detection model increases [46].

5.2 Disadvantage of using low-resolution image

Downsampling image leads to information loss and affects the accuracy of the model [41].

6 Discussion

3D data could be processed using 3D object detection models, or they could be sliced and fed to 2D object detection techniques for prediction. The 3D object detection technique, whether it is one stage approach or two-stage approach has these limitations:

- A major drawback of these 3D algorithms is that they are computationally expensive [49].
- 3D data are quite difficult to collect, annotate and store and hence there are fewer datasets for deep learning models. [47, 8]

In [44] to address the issue of more memory requirement and higher computational complexity Song et al. used low-resolution image as input for the Region Proposal Network(RPN). So, low-resolution image could be used to solve the problem of computational complexity and high memory requirement of 2D as well as 3D object detection networks. Reduction in resolution would also lead to reduction in accuracy. But the papers discussed in the sections one stage approach 2 and two-stage approach 3 don't provide the details about the impact of resolution on accuracy and speed of the models.

References

- [1] Adobe. <https://helpx.adobe.com/photoshop/kb/advanced-cropping-resizing-resampling-photoshop.html>. Accessed on 22 October 2020.
- [2] NIH Clinical Center. <https://nihcc.app.box.com/v/DeepLesion/folder/50715173939>. Accessed on 22 October 2020.
- [3] Hao Chen, Qi Dou, Lequan Yu, and Pheng-Ann Heng. Voxresnet: Deep voxelwise residual networks for volumetric brain segmentation. In *arXiv:1608.05895v1 [cs.CV]*, 2016.
- [4] Yilun Chen, Shu Liu, Xiaoyong Shen, and Jiaya Jia. Fast point r-cnn. 2019.
- [5] covid-19-chest-xray-lung-bounding-boxes dataset. <https://github.com/GeneralBlockchain/covid-19-chest-xray-lung-bounding-boxes-dataset>. Accessed on 22 October 2020.
- [6] Jifeng Dai, Yi Li, Kaiming He, and Jian Sun. R-fcn: Object detection via region-based fully convolutional networks. In *Advances in neural information processing systems*, pages 379–387, 2016.
- [7] M. Engelcke, D. Rao, D. Z. Wang, C. H. Tong, and I. Posner. Vote3deep: Fast object detection in 3d point clouds using efficient convolutional neural networks. 2017.
- [8] Di Feng, Xiao Wei, Lars Rosenbaum, Atsuto Maki, and Klaus Dietmayer. Deep active learning for efficient training of a lidar 3d object detector. In *arXiv*, 2019.
- [9] Andreas Geiger, Philip Lenz, and Raquel Urtasun. Are we ready for autonomous driving? the kitti vision benchmark suite. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2012.
- [10] Ross Girshick. Fast r-cnn. In *In Proceedings of the IEEE international conference on computer vision*, pages 1440–1448, 2015.

- [11] Yulan Guo, Hanyun Wang, Qingyong Hu, Hao Liu, Li Liu, and Mohammed Bennamoun. Deep learning for 3d point clouds: A survey. In *arXiv:1912.12033v2 [cs.CV]*, 2020.
- [12] Chenhong He, Hui Zeng, Jianqiang Huang, Xian-Sheng Hua, and Lei Zhang. Structure aware single-stage 3d object detection from point cloud. In *CVPR₂₀₂₀*, 2020.
- [13] Kaiming He, Georgia Gkioxari, Piotr Dollar, and Ross Girshick. Mask r-cnn. In *Facebook AI Research (FAIR)*, 2018.
- [14] Kaiming He, Georgia Gkioxari, Piotr Dollar, and Ross Girshick. Mask r-cnn. In *arxiv.org*, 2018.
- [15] Licheng Jiao, Fan Zhang, Fang Liu, Shuyuan Yang, Lingling Li, Zhixi Feng, and Rong Qu. A survey of deep learning-based object detection. In *arXiv:1907.09408v2 [cs.CV]*, 2019.
- [16] Kaggle. <https://www.kaggle.com/andrewmvd/covid19-ct-scans>. Accessed on 22 October 2020.
- [17] Naji Khosravan and Ulas Bagci. S4nd: Single-shot single-scale lung nodule detection. In *arXiv:1805.02279v2 [cs.CV]*, 2018.
- [18] Evi Kopelowitz¹ and Guy Englehard. Lung nodules detection and segmentation using 3d mask-rcnn. In *arXiv*, 2019.
- [19] Kingsley Kuan, Mathieu Ravaut, Gaurav Manek, Huiling Chen, Jie Lin, Babar Nazir, Cen Chen, Tse Chiang Howe, Zeng Zeng[†], and Vijay Chandrasekhar[†]. Deep learning for lung cancer detection: Tackling the kaggle data science bowl 2017 challenge. In *IEEE*, 2017.
- [20] Alex H Lang, Sourabh Vora, Holger Caesar, Lubing Zhou, Jiong Yang, and Oscar Beijbom. Pointpillars: Fast encoders for object detection from point clouds. In *CVPR*, 2019.
- [21] Andréanne Lemay. Kidney recognition in ct using yolov3. In *In Advances in neural information processing systems, pages 91–99*, 2019.
- [22] Bo Li. 3d fully convolutional network for vehicle detection in point cloud. In *IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, 2017.

- [23] Tsung-Yi Lin, Piotr Doll'ar, Ross Girshick, Kaiming He, Bharath Hariharan, and Serge Belongie. Feature pyramid networks for object detection. In *In Advances in neural information processing systems*, pages 91–99, 2015.
- [24] Tsung-Yi Lin, Priya Goyal, Ross Girshick, Kaiming He, and Piotr Doll'ar. Focal loss for dense object detection. In *In Proceedings of the IEEE international conference on computer vision*, pages 2980–2988, 2017.
- [25] Li Liu, Wanli Ouyang, Xiaogang Wang, Paul Fieguth, Jie Chen, Xinwang Liu, and Matti Pietikäinen. Deep learning for generic object detection: A survey. In *International Journal of Computer Vision (2020)* 128:261–318, 2019.
- [26] Wei Liu, Dragomir Anguelov, Dumitru Erhan, Christian Szegedy, Scott Reed, Cheng-Yang Fu, and Alexander C Berg. Ssd: Single shot multi-box detector. In *European conference on computer vision*, pages 21–37. Springer, 2016.
- [27] Jerubbaal John Luke, Rajkumar Joseph, and Mahesh Balaji. Impact of image size on accuracy and generalization of convolutional neural networks. In *IJRAR*, 2019.
- [28] microscope.org. <https://microscope-microscope.org/microscope-info/image-resolution/>. Accessed on 22 October 2020.
- [29] NIH. <https://www.nibib.nih.gov/science-education/science-topics/computed-tomography-ct>. Accessed on 22 October 2020.
- [30] NIH. <https://www.nibib.nih.gov/science-education/science-topics/magnetic-resonance-imaging-mri>. Accessed on 22 October 2020.
- [31] Radiological Society of North America. <https://www.kaggle.com/c/rsna-pneumonia-detection-challenge>. Accessed on 22 October 2020.

- [32] Charles R. Qi, Or Litany, Kaiming He, and Leonidas J. Guibas. Deep hough voting for 3d object detection in point clouds. 2019.
- [33] Charles R Qi, Wei Liu, Chenxia Wu, Hao Su, , and Leonidas J Guibas. Frustum pointnets for 3d object detection from rgb-d data. 2018.
- [34] AZO quantum. <https://www.azoquantum.com/News.aspx?newsID=7018>. Accessed on 22 October 2020.
- [35] Joseph Redmon, Santosh Divvala, Ross Girshick, and Ali Farhadi. You only look once: Unified, real-time object detection. In *In Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 779–788, 2016.
- [36] J. Redmon and A. Farhadi. Yolo9000: Better, faster, stronger. In *In CVPR*, 2017.
- [37] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. Faster r-cnn: Towards real-time object detection with region proposal networks. 2015.
- [38] Shaoqing Ren, Ross Girshick Kaiming He, and Jian Sun. Faster r-cnn: Towards real-time object detection with region proposal networks. In *In Advances in neural information processing systems*, pages 91–99, 2015.
- [39] Shaoshuai Shi, Xiaogang Wang, and Hongsheng Li. Pointnet: 3d object proposal generation and detection from point cloud. In *In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 770–779, 2019.
- [40] Shaoshuai Shi, Zhe Wang, Xiaogang Wang, and Hongsheng Li. Part-a2net: 3d part-aware and aggregation neural network for object detection from point cloud. 2019.
- [41] Connor Shorten and Taghi M. Khoshgoftaar. A survey on image data augmentation for deep learning. In *Springer*, 2019.
- [42] M. Simony, S. Milzy, K. Amendey, and H.-M. Gross. Complex-yolo: an euler-region-proposal for real-time 3d object detection on point clouds. In *In Proceedings of the European Conference on Computer Vision (ECCV)*, pages 0–0, 2018.

- [43] S. Song, S. Lichtenberg, and J. Xiao. Sun rgb-d: A rgb-d scene understanding benchmark suite. In *Proceedings of 28th IEEE Conference on Computer Vision and Pattern Recognition (CVPR2015)*, 2015.
- [44] Shuran Song and Jianxiong Xiao. Deep sliding shapes for amodal 3d object detection in rgb-d images. 2016.
- [45] Tang, Y., Yan, K., Tang, Y., Liu, J., Xiao, J., Summers, and R.M. Uldor: A universal lesion detector for ct scans with pseudo masks and hard negative example mining. In *arXiv:1901.06359*, 2019.
- [46] towards data science. <https://towardsdatascience.com/no-gpu-for-your-production-server-a20616bb04bd>. Accessed on 22 October 2020.
- [47] Qian Wang, Neelanjan Bhowmik, and Toby P. Breckon. Multi-class 3d object detection within volumetric 3d computed tomography baggage security screening imagery. In *arXiv*, 2020.
- [48] Yan Yan, Yuxing Mao, and Bo Li. Second: Sparsely embedded convolutional detection. In *MDPI*, 2018.
- [49] Anqi Yang, Aswin Sankaranarayanan, Srinivasa Narasimhan, David Held, and Jen-Hao Chang. 3d object detection from ct scans using a slice-and-fuse approach. In *Carnegie Mellon University*, 2019.
- [50] Bin Yang, Wenjie Luo, and Raquel Urtasun. Pixor: Realtime 3d object detection from point clouds. In *In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pages 7652–7660*, 2018.
- [51] Zetong Yang, Yanan Sun, Shu Liu, Xiaoyong Shen, and Jiaya Jia. Std: sparse-to-dense 3d object detector for point cloud. In *ICCV*, 2019.
- [52] Zhengdong Zhang, Shuai Li, Ziyang Wang, and Yun Lu. A novel and efficient tumor detection framework for pancreatic cancer via ct images. In *arXiv*, 2020.
- [53] B. Zhao, J. Soraghan, G. Di-caterina, L. Petropoulakis, D. Grose, and T. Doshi. Automatic 3d segmentation of mri data for detection of head and neck cancerous lymph nodes. In *2018 Signal Processing: Algorithms, Architectures, Arrangements, and Applications (SPA)*, pages 298–303, 2018.

- [54] Yin Zhou and Oncel Tuzel. Voxelnet: End-to-end learning for point cloud based 3d object detection. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2017.
- [55] Wentao Zhu, Chaochun Liu, Wei Fan, and Xiaohui Xie¹. Deeplung: Deep 3d dual path nets for automated pulmonary nodule detection and classification. In *arXiv:1801.09555v1 [cs.CV]*, 2018.