



Real-Time 3D Object Reconstruction and Distance Estimation from 2D Images using OpenCV and Unity

Kotha Abhishek

Department of Artificial Intelligence & Machine Learning, GMR Institute of Technology, Rajam, Andhra Pradesh, India
kothaabhishek125@gmail.com

ABSTRACT—

This paper introduces a methodology that combines OpenCV, Unity, and Generative Adversarial Networks (GANs) to achieve real-time 3D object reconstruction and precise distance estimation from 2D images. Leveraging OpenCV's robust computer vision, Unity's interactive platform, and GANs' advanced 3D modeling capabilities, the approach transforms static 2D images into dynamic 3D models, enriched with live distance data. Through intricate feature extraction and depth analysis, the system generates interactive 3D representations within Unity, significantly narrowing the gap between 2D images and immersive 3D models. The immediate and accurate distance estimations offered by this system hold great promise for applications in augmented reality, virtual simulations, and interactive design. Beyond the realm of virtual experiences, the generated 3D models find practical applications in robotics. The precise distance estimations derived from reconstructed objects serve as invaluable inputs for robotic systems, enhancing their navigation and interaction capabilities in the real world. Additionally, this paper envisions a positive impact on the visually impaired community, as the real-time distance estimations can be utilized to create applications that enhance spatial awareness, fostering inclusivity and accessibility. Furthermore, the paper explores the potential for entertainment applications by applying skeletonization techniques to the generated 3D models. This opens up avenues for creating augmented reality (AR) games, providing users with immersive and interactive gaming experiences fuelled by realistic 3D models. In conclusion, the integration of OpenCV, Unity, and GANs in this approach not only promises advancements in spatial understanding but also unlocks diverse possibilities for practical applications, spanning robotics, accessibility, and engaging digital entertainment.

Keywords- : OpenCV, 2D to 3D reconstruction, Target region extraction, Deep Learning, Dimension Measurement, Unity, GAN.

1. INTRODUCTION

Nowadays, real-time perception has emerged as a pivotal advancement, revolutionizing the realms of robotics and self-driving cars. The integration of 3D model visualization has notably amplified the performance of motion-capable humanoid robots, ushering in a new era of precision and agility. Additionally, this technological stride holds immense promise for empowering the visually impaired community.

This paper delves into an innovative approach that leverages depth estimation to intricately pinpoint objects within images. The synergy of depth data and 3D models represents a concerted effort to achieve unparalleled accuracy in object localization. The integration of deep learning techniques plays a pivotal role in this process, aiding in the identification of priority objects and facilitating depth determination through the use of bounding boxes. A noteworthy aspect of this methodology lies in its capacity to not only ascertain object size but also accurately gauge the distance of these objects from the observer. In the dynamic realm of moving objects, our approach allows for the analysis of parameters such as speed, albeit with inherent challenges arising from variations in viewer angles and image quality.

An additional layer of innovation comes with the integration of Generative Adversarial Networks (GANs), acting as a transformative force within our framework. Through the application of adversarial training, GANs contribute significantly to enhancing the precision and realism of 3D models. This infusion of GANs into the real-time perception framework holds the promise of a revolutionary shift, further solidifying the impact of this technology on various applications.

In summary, the amalgamation of depth estimation, 3D models, deep learning, and GANs represents a holistic and groundbreaking approach to real-time perception. This not only paves the way for advancements in robotics and autonomous vehicles but also extends its reach to benefit the visually impaired community, promising a future where precision and accessibility coalesce in unprecedented ways.

Literature Survey

The paper mentions that there has been immense progress in neural scene representations, with various methods optimized from 2D multi-view images using differentiable neural rendering. Implicit representations, such as coordinate networks, employ large MLPs to model 3D scenes, offering advantages

in memory efficiency. Explicit representations, like voxels, are fast to query but memory-intensive, making them difficult to scale up for large-scale scenes. Hybrid representations combine the benefits of both explicit and implicit representations to strike a balance between efficiency and detail. SceneDreamer is an unconditional generative model for unbounded 3D scenes, learned from in-the-wild 2D image collections without 3D annotations [1].

The paper discusses the importance of visual perception in autonomous driving and simultaneous localization and mapping (SLAM). A deep learning method, Convolutional Neural Network (CNN), was used to detect the bounding boxes around the objects. The occultation problem has been overcome for the single object model construction. There is poor 3D visualization. It performs very quickly for closed area model development. It gives the better results on RMSE errors. However, poor performance was observed at high frame rates. The importance of 3D object detection and localization in mobile robot navigation, and the use of anchors for improving the performance of object detection and distance estimation [2].

You Only Look Once (YOLO) method was used for object identification and Region-Based Convolutional Neural Networks (R-CNNs) for model performance and analysis. Dimensions were measured using OpenCV. Canny edge detector was used to improve dimension detection analysis. This approach was very fast. It showed less accuracy in dimension measurement for large objects. Works with just a webcam with some high quality. The paper emphasizes the use of Python programming language for consistency, simplicity, and access to excellent libraries and frameworks in ML and AI-based papers. It highlights the benefits of using OpenCV functions to determine an object's length, breadth, and volume [3].

The paper mentions the use of stochastic approaches to avoid high nonlinearity in motion estimation for a chain of joints in the human body. Converts the human into 3D models with the same pose. Adjusts the image into the correct position. Works quickly for real-time videos and images. Does not require high hardware settings. Nonlinear iterative algorithm used. For the human detection CNN is used. Gives the high accuracy. The paper focuses on real-time applications and considers the heavy computational load as a critical drawback. The paper introduces rotation angle variables and unit vectors in the direction of the rotation axis to formulate the motion estimation process [4].

Constructed the 3D indoor room layout. Deep learning method CNN is used. Topology anchor point optimization (TAPO) is used. TAPO gives the more accuracy. Public data sets LSUN, Hedau and 3DGP are used for the model train. Gives the complete details like the area and length and breadth. Do not give the other objects in the room. Previous works have extensively investigated layout estimation due to its significance in scene understanding and object recognition. Floors typically serve as supporting surfaces for objects like chairs and tables [5].

The paper provides a comprehensive review of object detection techniques using deep learning, aiming to achieve high accuracy and real-time performance. It discusses the challenges faced by object detection systems that rely on other computer vision techniques, leading to slower and less efficient performance. The paper focuses on using a completely deep learning-based approach to solve object detection problems in an end-to-end fashion using wireless sensor networks. It also mentions the general applications and results of object detection using deep learning, highlighting the importance of object recognition in various industries, such as video surveillance and search engines. The paper emphasizes the need for robust feature descriptors to accurately represent different types of objects, mentioning SIFT, HOG, and hair-like features as commonly used representations [6].

Layered depth images (LDI) representation is used for multi-view synthesis to address occlusion and hidden information in a 3D scene. Vincent Janett, et al., proposed an improved virtual synthesis based on object levels distinction using region growing segmentation technique. The paper introduces a cost-effective approach using a single view with multi-colour filter aperture (MCA) and multi-plane representation for 3D image generation. A 2.1D sketch is used as a semantic segmentation technique to determine the number of objects and identify occlusion in the 3D scene. Experimental validations are conducted to validate the proposed approach with depth gaps ranging from 0.5cm to 10.5cm [7].

3D target detection has been a research hotspot in recent years. Methods for 3D target detection can be divided into feature descriptor-based, template-based matching, and learning-based approaches. Feature descriptor-based methods focus on RGB images and use feature point matching algorithms to find corresponding 2D pixel points in the template image. Representative methods in this category include SIFT, SURF, FAST, and ORB. Template-based matching methods require multi-view imaging of the target object and extract contour, surface curvature, and edge information. Examples of methods in this category include Linemod, SSD-6D, PoseCNN, Recovery 6d, and Robust 6d. Learning-based methods use deep learning techniques to detect 3D targets in indoor scenes, as demonstrated in the current paper [8].

The first model proposed for 2D to 3D conversion was Deep3d, which used a neural network for disparity map prediction and generated the right view of an image from the left view. Deep3d had limitations such as a fixed image size and non-differentiability. To address these, a distribution over a range of disparity values was predicted at each pixel, and the right image was obtained by computing the expected value over the range of disparity values. Another approach proposed a fully convolutional deep neural network inspired by the supervised DispNet architecture. This approach treated depth estimation as an image reconstruction problem and used a left-right consistency check to improve the quality of synthesized depth images. The goal of this approach was to predict the per-pixel scene depth given a single image, without requiring ground truth depth. It aimed to find the dense correspondence field that would enable the reconstruction of the right image from the left image [9].

The paper provides a comprehensive literature review of target detection and discusses works closely related to it. Various object detection methods, including one-stage and two-stage detectors, are systematically summarized in the paper. The datasets and evaluation criteria used in object detection are introduced. The development of object detection technology is reviewed in the paper. The paper also discusses the main research directions in the future based on the understanding of the current development of target detection [10].

The paper discusses MiddleVR, which is middleware for the Virtual Reality environment, handling head tracking, controls, and camera creation in the game engine. The paper also mentions the use of Unity as the game engine in the project, along with the use of MiddleVR. It is noted that the Unity

program in the project could only run at around 17 FPS, which was attributed to the different materials on each cube. To overcome this issue, the paper suggests using Unity's Prefab feature, which allows for the creation of objects consisting of multiple individual objects. Additionally, the paper mentions the use of raw DICOM data, which can be converted into .tiff image files for use in Unity. Overall, the paper focuses on the use of MiddleVR and Unity in a project, addressing performance issues and the utilization of DICOM data [11].

The paper discusses the application of OpenCV for feature extraction from 2D engineering drawings and the reconstruction of 3D CAD models. It highlights the availability of legacy designs in mechanical, aerospace, and civil engineering as drawings rather than software-generated CAD models. The proposed methods aim to convert such drawings into 3D CAD and BIM models using camera capture or scanned drawing data. The paper acknowledges that the method currently does not account for hidden lines in CAD drawings. The work leverages the latest developments in computer vision to develop an algorithm for automated conversion of 2D engineering drawings to 3D CAD models. SCAD geometries are generated from 2D geometries in this process. The methodology is developed as a python Jupyter-notebook, which is easily accessible and can be run locally. The growth of 3D modeling software has facilitated rapid 3D prototyping using CAD, and this work contributes to the automated conversion of 2D drawings to 3D models [12].

Significant research has been conducted on 3-D city modelling using aerial imagery, airborne LiDAR, ground-based or vehicle-borne sensing techniques, and combinations thereof. Early research focused on automatic extraction of building data from aerial imagery due to its high spatial resolution. Recent studies have relied on airborne LiDAR-based rooftop modelling, mobile LiDAR modelling of building facades, and hybrid approaches for large-scale urban modelling. Approaches to creating 3-D building models from aerial images on a large scale are categorized into parametric shapes, segmentation, and digital surface models. Various methods have been developed for generating digital terrain models (DTMs) that represent the elevations of 3-D urban models. Laser scanners can generate more accurate 3-D geometric primitives than image-based representations, but airborne LiDAR data have a low point density on building facades [13].

There are various 3D reconstruction techniques, including explicit measurements with laser or radar sensors, using multiple images, and using video sequences. Most vision-based approaches have focused on stereovision and algorithms that require multiple images, such as optical flow, structure from motion, and depth from defocus. Some algorithms reconstruct the 3D shape of known objects from images and laser data. Structured lighting offers another method for depth reconstruction. There are algorithms that can perform depth reconstruction from single images in specific settings, such as surface reconstruction for known objects like hands and faces. Methods like shape from shading and shape from texture generally assume uniform colour and texture, which may not perform well on complex, unconstrained, highly textured images [14].

A stereovision based method was used to estimate the distance with high accuracy in real time. The triangulation method was used to calculate the distance. The method achieved an average accuracy of 95%. The computation speed was also very high. Kernelized correlation filter (KCF) was used for object tracking. Worked for the live video also. Do not work for the ordinary cameras. The experimental results show that the proposed method is suitable for a baseline of 151mm and can estimate distances accurately from 88cm to 300cm. Each pixel's disparity is further refined using a simple but effective approach [15].

2. METHODOLOGY

2.1 Scene Dreamer

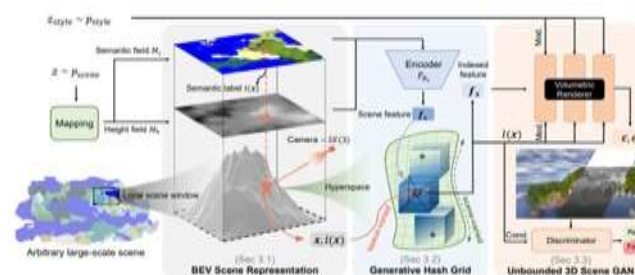


Fig -1: working of scene Dreamer model

The fig-1 shows that generalised steps for how the scene Dreamer is going work.

Data Collection:

In this paper they can download the large amount of real-life wild images from the google. After they can preprocess them by intensity of the light of the picture. Finally, they can train their model with the accurate images of the wild.

Scene Dreamer

In this first it takes the 2D image as the input and identify the small portion which is realistic to represent the 3D model and after the selection of that portion then it going to generate the height representation model of that image by using the depth estimation of the each pixel in the image and generate the BEV representation. After that it going to find the semantic field which contains the information of the natural items which in the image like lakes,

mountains, rivers etc. Now the generative hash grid is going to generate and we can perform the operation on that grid map. By using the GAN model this can provide the realistic 3D model by taking the grid map and finally give the model. The GAN can work like the generator and the detector of the faults. First it generate the model and it is only going to rectify the in it like the updating of the weights in the deep learning.

The generated model can also used for the games to if we can add the skeleton for the generated model.

While using this model can only generated the bird eye view only but we cannot zoom it and going to observer the specific position clearly. The model developing also the data set will need is very high and the accurate images. The other models also there but this can generate the real world distance images by using the single RGB image. But the other applications will need the specific Chamara like the stereo and the other specific one.

2.2: YOLO

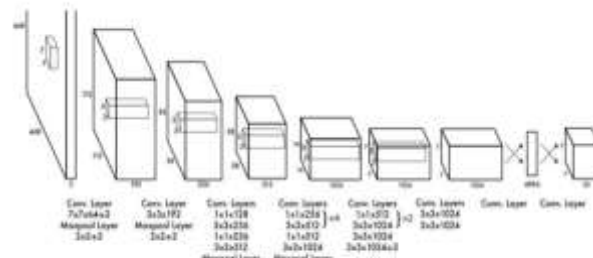


Fig 4. Convolution layers

Fig. – 2: Convolution layers

Data Collection:

In this paper they use the Image Net data set for training and the testing their model. This data set consists of the large no of the label images data. From this we can just prefer what images we need to train their model.

Working:

YOLO is a single-shot object detection algorithm, meaning it detects objects in a single pass through the neural network. It divides the input image into a grid and predicts bounding boxes and class probabilities for objects within each grid cell, enabling it to detect multiple objects at once.

YOLO predicts bounding boxes for possible object locations, along with class probabilities indicating the likelihood of an object belonging to a specific category. YOLO's single forward pass through the network makes it significantly faster than two-stage detectors, making it suitable for real-time object detection tasks. YOLO is widely applied in applications such as image and video analysis, robotics, autonomous vehicles, and more, due to its efficiency and accuracy in detecting objects. The detection of the objects is done by the CNN. The bounding boxes will come in the image. Now the SVM is going to classify the images based on the label. So that two algorithms going to classify the images effectively. The main advantage of the YOLO is that for the one time see then it going to identify if the same is come again. By this the training time is going to less. And the it going to find the labels with one pass of the image. So it is also going to use in the daily life use also like the video editing, and the live object detection.

2.3 Polygonal Model

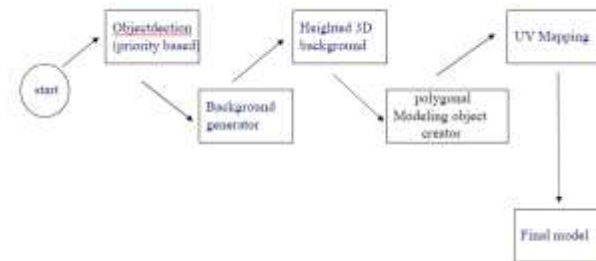


Fig. – 4: workflow of the polygonal model.

Data Collection:

In this paper they can take the daily life objects to convert into the 3D models. So they can gather images from different datasets and merger them into the single data set for there use.

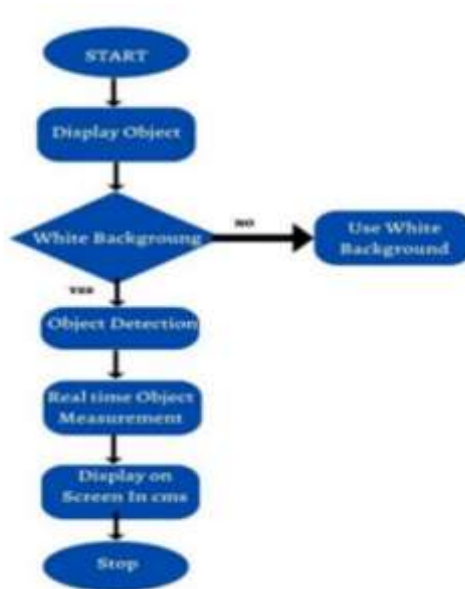
Working:

A 2D image is going to take and then remove the noise in the image (less priority objects in the image). After that we going to find the labels of the objects present in the image . And then a 2D background is going to generate and blend that with the 3D height model using the Depth estimation of the region. Now the label objects is going to generate according to the required size by using polygonal model now some new objects is create using this data . Now UV mapping is going to happen after that normalize the model for the better realistic effect and add the tuning the model for the more fine details. If we want the we also add the skeleton for the model for using the gaming purpose. The main advantage of the polynomial model is that we can represent any object into the 3D models. By getting this 3D models of the objects they place them in the required position by using the unity. As per the theory that every object can represent in the different polygons. So for the difficult daily life objects like the pen, and the ball like small objects can also going to represent in the polygonal method.

The disadvantage of the polynomial model is that it do not create the infinite land scraps into the

3D model. It only converts the objects into models for the 3D conversion is already object detection the ML models can use we generate the 3D model of the object base

2.4 Open CV modules:



the 3D models. The pretrained available in the market. For the can get the labels and then we can on the size of the image

Fig. – 5: work flow of the object size estimation.

In this paper the basic code is used for the detection of the real time object measurement.

First objects detected by the yolo can be taken by this algorithm. Then check for the closed contours. Then by the simple code we can find the Area , width, height and the center of the Object and going to display. This approach can only use the basic functions in the OpenCV for this work.

The use of this method is for the ratio of the space it is going to occupy in the 3D model. Based on its volume we can put this in the 3D model accurately. The contour finding we can continue method for the getting closed object in the image. The region growing or other things also use full for the clear object detection in the bounding box. For this also we can use the basic OpenCV basic modules for the detection. There are many OpenCV function are available for the better contour finding in the image.

```

Describing the size of objects in the picture
cv2.putText(orig, "l: {:.1f}cm".format(lebar_pixel/25.5),
            int(trbex + 10), int(trbpy), cv2.FONT_HERSHEY_SIMPLEX, 0.7, (0,0,255), 2)
cv2.putText(orig, "P: {:.1f}cm".format(panjang_pixel/25.5), (int(tltx - 15),
            int(tlty - 10)), cv2.FONT_HERSHEY_SIMPLEX, 0.7, (0,0,255), 2)
cv2.putText(orig, str(areal), (int(x), int(y)))
cv2.FONT_HERSHEY_SIMPLEX, 0.6, (0,0,0), 2)
hitung_objek+=1

#Displays Number of Detected Objects
cv2.putText(orig, "OOP-EL: {}".format(hitung_objek), (10,50),
            cv2.FONT_HERSHEY_SIMPLEX, 1, (0,0,255), 2, cv2.LINE_AA)
cv2.imshow('Camera', orig)
  
```

Fig. – 6: Example code for object size estimation.

2.5 Convolutional Neural Network:

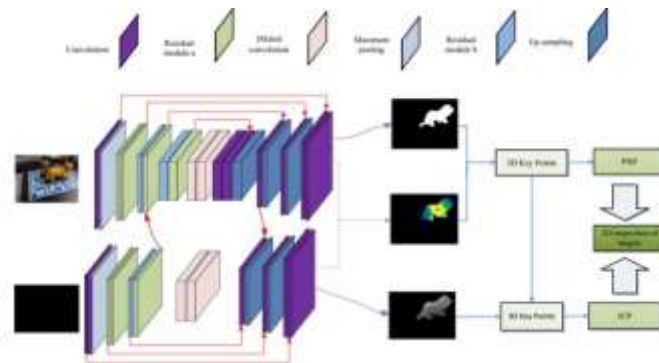


Fig. – 7: : working of the convolutional neural network.

The paper proposes a deep learning-based approach for 3D target detection in complex indoor scenes. The authors construct an indoor 3D target detection dataset using Aruco markers and establish a pixel-by-pixel key point voting network for joint semantic segmentation of RGB images.

They propose a new key point assumption strategy and extend key point detection to three dimensions using depth images. The paper demonstrates the generalization of the method and achieves 3D target detection in indoor scenes based on RGB images and RGB-D images.

The evaluation metrics and visualization of the model are analyzed and compared, including testing and visualization under validation set, truncated validation set, and unlabeled scenes. This method can use in the two types the fast retrieving of the objects from the image. This method is more efficient for the identifying the specific pattern in the image. Like if the racks need to find means first we need to prepare the filter bases on the pattern of the racks. Then we can use that filter on the image to find that pattern is there in the image or not. This method is not suitable for the identifying the complex objects. For them we can use the basic machine learning algorithms for finding them.

This is the complete over view of the methodology in this paper.

Comparison Table:

	Title	Year	Objectives	Limitations	Advantages
Reference 1	Un bounded 3D Scene Generation from 2D Image Collections.	2023	Making the 3D models for unbounded images. Improve the 3D vision of using opencv.	Inaccurate for the lot of the Noise image.	Compleate 3D models in the bird eye view. We can generate them with single RBG image.
Reference 2	Ancher Distance for the 3D Multi-Object Diistance Estimation From 2D single shot.	2021	By converting the images to 3D models can improve the auto break feature better. Increase the accuracy of the self driving cars	In this they only done for the cars identification only. They do not represent 3D models in the realistic way	For the cars they can overcome the object occlusion problem. Better real time distance estimation for the objects
Reference 3	Real Time Object distance and Dimension Measurement using Deep Learning and OpenCV	2023	Find the Object dimension in cm using the OpenCV.	Inaccurate for the Large Objects. Need the distance of the image and the camara.	Done using the basic function in the OpenCV.
Reference 4	2D-to-3D Visual Human Motion Converting System for Home Option Motion Capture Tool and 3-D smart TV	2015	Converting the 2D human phose into the 3D model.	This model is for the home only. It detect only the single one.	It do not need any other things only the model installatied tv. Increase the VR games effectively.
Reference 5	3D room Layout Estimation From a Single RGB image	2020	Construct the 3D model room from the 2D RGB Image. Estimate the area of the room.	Scale ambiguity for the detection the floor. Camara need to cover all the walls and ceil.	They use the CNN for the construction of the room. Better version of the feature civil.
Reference 6	Object Dection using Deep Learning.	2021	Detecting the objects in the image.	No of the objects in the image. Needed the clear image	Using the deep learning models. Yolo for the fast retrieving

				with less noise.	the objects.
Reference 7	Accuracy in Depth Recovery and 3D image synthesis Fromm Single Image Using Multi-Color Filter Aperture and Shallow Depth of field.	2021	Estimate the Depth of the objects in the image. Increase the precision of the 3D model.	Need the specila RGB filterd Chamara. Need the clear back ground.	For the particular object the 3D model is the accurate.
Reference 8	Deep learning based 3D target detection for indoor scenes.	2023	Detects the specific object in the complete image. Converts the detected object to 3D Object.	Detect only the particular object. Create only 3D model for the particular object.	Give accurate detection of the particular object. Better for thehome appliances finding.
Reference 9	2D to 3D image conversion Using machine Learning Approach.	2019	Converting of the normal images into stereo pair images. Converting images/movies into the 3D images/movies.	Blure images are hard to convert into the 3D models.	We do need the stereo cameras. Normal images into the 3D images.
Reference 10	Deep Learning for Object Dection.	2021	Detect the object in the given images, which are having the area greater than some particular value.	Small object detection is difficult. Large data set and time required to train the model.	Also reorganize the objects if the image slightly having the blur. Overcome the occlusion.
Reference 11	Visualization of MRI scan images as a 3D object within Unity.	2019	Visualizing the 2D MRI scans into the 3D models.	In this approach of this paper the back site cubes will not visible.	Powet full unity is used for the representation the 3D models. Improve the medical field.
Reference 12	Automated 3D solid reconstruction from 2D CAD using OpenCV	2021	Constructing the 3D model using the 2D CAD models.	Only the CAD images are convertible.	Can achive better visualisation of the objects in the hardware.
Reference 13	Improving accuracy of the automated 3D building models for smart cities.	2017	Create 3D models of the cities. Improve the horizontal and vertical accuracy of the adjusted 3D buildings.	More time for generating the model. High quality cameras are required.	Improves the navigation system. Accurate 3D city models enabel better decision making in urban planning and management.
Reference 14	3D Depth Reconstruction from single still Images.	2019	The paper aims to address the challenging task of estimating 3D depth from a single still image using a supervised learning approach.	Gives the low quality depth map for the long distance images.	Increase the quality of the depth estimation for the outdoor features.
Reference 15	Moving Object Distance Estimation method Bases on Target Extraction with a Stereo camera.	2019	Reorganize the moving objects with accurate.	Must use the stereo camera. Focus on the single Object.	Less computation time. High accuracy for the moving objects also.

3. GRAPHICAL REPRESENTATION



Fig.-8: Graphical Representation of different authors with respective years

In the above graph, X- axis represents the authors and Y- axis represents the published year. Algorithms that are framed in the picture is the best in the respective papers.

4. RESULTS

The below table incorporates with the results of 5 different reference papers. It includes the datasets, methods used and performance metrics of the respective papers. Most of the authors preferred Scene Dreamer model because it is getting the better result than any model in the present tech, making the effective for creating the real world 3D models. And also worked for the data which is bounder less images also.

Reference Number	Method used	Performance Metrics
[1]	Scene Dreamer	KID = 4.52 Fid = 76.73
[16]	Gan craft	KID = 4.51 Fid = 79.99
[6]	Yolo	Accuracy = 95% MAE:1.23
[8]	R-cnn	Accuracy= 93%
[3]	Basic operations	Accuracy = 91% RMSE = 0.52

Accuracy: Accuracy is a common metric used to evaluate the performance of classification models. It represents the ratio of correctly predicted instances to the total number of instances in the dataset.

Root Mean Square Error (RMSE): Root Mean Square Error (RMSE) is a metric used to measure the average magnitude of the errors between predicted values and actual values. It is often used in the context of regression analysis and machine learning to evaluate the performance of a predictive model.

Mean Absolute Error (MAE): Mean Absolute Error (MAE) is another metric used to measure the average magnitude of errors between predicted values and actual values. Like RMSE, MAE is often used in regression analysis and machine learning to assess the performance of a predictive model. However, MAE differs from RMSE in the way it calculates the error.

5. CONCLUSION

The 'Scene Dreamer' algorithm has proven to be a leading choice for creating outdoor 3D models from single-shot 2D images, demonstrating superior performance even on challenging unbounded images. While the running time for generation is relatively longer compared to other models, the algorithm excels in capturing fine details of all objects in the image. In this model it only generate only the bird eye view with the normal rgb image. But the all other model need the specific feature of the image like the stereo camera or the other specific needs.

For object detection, YOLO (You Only Look Once) stands out as the optimal model. It is renowned for its efficiency and real-time capabilities. However, the selection of the 'best' model depends on specific task requirements, data characteristics, and various factors. The yolo model can predict the object bounding boxes in the one pass of the image. And the algorithm is quite difficult to understand. For the more accuracy this yolo using the two algorithms one is for the object detection and other is for the object classification.

In conclusion, the 'Scene Dreamer' algorithm offers remarkable performance in outdoor 3D model generation, particularly excelling in challenging scenarios. Meanwhile, YOLO remains a top choice for efficient object detection. The generation of the model can also need to fix bases on the our need. While if we need the indoor scenes then we can give high priority for the indoor objects and neglect the outdoor small things like the bushes and the trees. So this method quiz efficient for the all kind of the uses in our daily needs. The decision on the most suitable model should be based on the unique needs and considerations of the task at hand.

REFERENCES

- [1] Chen, Z., Wang, G., & Liu, Z. (2023). Scenedreamer: Unbounded 3d scene generation from 2d image collections. arXiv preprint arXiv:2302.01330.
- [2] Yu, H., & Oh, J. (2021). Anchor distance for 3d multi-object distance estimation from 2d single shot. *IEEE Robotics and Automation Letters*, 6(2), 3405-3412.
- [3] Basavaraj, M. U., & Raghuram, H. (2023, February). Real Time Object Distance and Dimension Measurement using Deep Learning and OpenCV. In *2023 Third International Conference on Artificial Intelligence and SDmart Energy (ICAIS)* (pp. 929-932). IEEE.
- [4] Han, Y. (2014). 2D-to-3D visual human motion converting system for home optical motion capture tool and 3-D smart TV. *IEEE Systems Journal*, 9(1), 131-140.
- [5] Yan, C., Shao, B., Zhao, H., Ning, R., Zhang, Y., & Xu, F. (2020). 3D room layout estimation from a single RGB image. *IEEE Transactions on Multimedia*, 22(11), 3014-3024.
- [6] Pal, S. K., Pramanik, A., Maiti, J., & Mitra, P. (2021). Deep learning in multi-object detection and tracking: state of the art. *Applied Intelligence*, 51, 6400-6429.
- [7] Deshpande, R. R., Bhatt, M. R., & Madhavi, C. R. (2021). Accuracy in depth recovery and 3D image synthesis from single image using multi-color filter aperture and shallow depth of field. *IEEE Access*, 9, 123528-123540.
- [8] Liu, Y., Jiang, D., Xu, C., Sun, Y., Jiang, G., Tao, B., ... & Yun, J. (2023). Deep learning based 3D target detection for indoor scenes. *Applied Intelligence*, 53(9), 10218-10231.
- [9] Fu, J., Yang, Y., Singhrao, K., Ruan, D., Chu, F. I., Low, D. A., & Lewis, J. H. (2019). Deep learning approaches using 2D and 3D convolutional neural networks for generating male pelvic synthetic computed tomography from magnetic resonance imaging. *Medical physics*, 46(9), 3788-3798.
- [10] Xiao, Y., Tian, Z., Yu, J., Zhang, Y., Liu, S., Du, S., & Lan, X. (2020). A review of object detection based on deep learning. *Multimedia Tools and Applications*, 79, 23729-23791.
- [11] DE GRAAF, L. I. Z. A. Visualization of MRI scan images as a 3D object within Unity.
- [12] Harish, A. B., & Prasad, A. R. (2021). Automated 3D solid reconstruction from 2D CAD using OpenCV. arXiv preprint arXiv:2101.04248.
- [13] Lim, B., & Zohren, S. (2021). Time-series forecasting with deep learning: a survey. *Philosophical Transactions of the Royal Society A*, 379(2194), 20200209.
- [14] Yang, B., & Lee, J. (2019). Improving accuracy of automated 3-D building models for smart cities. *International journal of digital earth*, 12(2), 209-227
- [15] Zhang, J., Chen, J., Lin, Q., & Cheng, L. (2019, July). Moving object distance estimation method based on target extraction with a stereo camera. In *2019 IEEE 4th International Conference on Image, Vision and Computing (ICIVC)* (pp. 572-577). IEEE.
- [16] Hao, Z., Mallya, A., Belongie, S., & Liu, M. Y. (2021). Gancraft: Unsupervised 3d neural rendering of minecraft worlds. In *Proceedings of the IEEE/CVF International Conference on Computer Vision* (pp. 14072-14082).