A

# COMPUTER VISION PROJECT REPORT

on

# Vision sense Depth Estimation

Submitted by:

## G Sai Nikhil(210214)
## D Veera Harsha Vardhan Reddy(210258)
## N Karthik Raja(210371)

under mentorship of

## Dr.Sukhandeep Kaur

Assistant Professor

Department of Computer Science Engineering
School of Engineering and Technology
BML MUNJAL UNIVERSITY, GURUGRAM (INDIA)

Dec 2024

# CANDIDATE'S DECLARATION

I hereby certify that the work on the project entitled, "**Vision Sense Depth Estimation**", in partial fulfilment of requirements for the award of Degree of **Bachelor of Technology** in School of Engineering and Technology at BML Munjal University, having University Roll No's : 210371, 210214 and 210258 is an authentic record of our own work carried out during a period from Aug 2024 to December 2024 under the supervision of Dr.Sukhandeep Kaur

(**Karthik Raja)**
(**Sai Nikhil)**
(**Harsha Vardhan Reddy**)

# SUPERVISOR'S DECLARATION

This is to certify that the above statement made by the candidate is correct to the best of my knowledge.

**Faculty Supervisor Name :** Dr.Sukhandeep Kaur
**Signature :**

**TABLE OF CONTENTS**

## Abstract

Depth estimation is a critical task in computer vision, enabling the extraction of spatial information from images for applications such as robotics, autonomous systems, and augmented reality. This study investigates monocular and binocular depth estimation using distinct methodologies tailored to each approach. Monocular depth estimation focuses on indoor environments, utilizing the NYU Depth V2 dataset and three deep learning models Deeplabv3, U-Net, and DenseNet. Preprocessing techniques such as edge detection, Gaussian noise application, and blurring were employed to analyse their impact on model performance. Binocular depth estimation, performed with the KITTI dataset, relied on traditional computer vision methods, including block matching, semi-global block matching, and MiDAS-Small, without the use of deep learning models. Preprocessing steps for binocular approaches included histogram equalization, edge detection, and Gaussian blurring to enhance stereo image quality. The report provides a comprehensive comparative analysis of these methods, highlighting the effectiveness, strengths, and challenges of both deep learning and conventional techniques in diverse depth estimation scenarios.

**Keywords:** Depth Estimation, Computer vision, Deep Learning, U-Net, DeeplabV3, Blurring

# 1 Introduction

Over the past 10 years, depth estimation has been a significant research area in computer vision due to its wide-ranging applications and the growing capabilities of computational techniques. Depth estimation (DE) is a traditional computer vision task that predicts depth from one or more two-dimensional (2D) images. DE estimates each pixel's depth in an image using offline-trained models. In machine perception, recognition of some functional factors such as the shape of a scene from an image and image independence from its appearance seems to be fundamental. DE has great potential for use in disparate applications, including grasping in robotics, robot-assisted surgery, computer graphics, and computational photography.

The DE task needs an RGB image and a depth image as output. The depth image often consists of data about the distance of the object in the image from the camera viewpoint. The computer-based DE approach has been under evaluation by various investigators worldwide, and the DE problem has been an exciting field of research. Most successful computer-based methods are employed by determining depth by applying stereo vision. With the progress of recent deep learning (DL) models, DE based on DL models has been able to demonstrate its remarkable efficiency in many applications. DE can be functionally classified into three divisions, including monocular depth estimation (MDE), binocular depth estimation (BDE), or multi-view depth estimation (MVDE).

Sometimes, the application of algorithms for calculating depth may create different challenges. For instance, the matching cost function utilized in the algorithm can generate false-positive signals, which eventuates in the creation of depth maps with low accuracy. Thus, the use of post-processing approaches (i.e., median filter, bilateral filter, and interpolation) is of great importance in stereo vision applications to delete noise and refine depth maps.

# 2 Literature Review

The paper[1] introduces a depth estimation method from single RGB images, which relies on depth-from-defocus cues rather than traditional multi-view approaches. Specifically, it proposes a unique Point Spread Function (PSF) convolutional layer that utilizes location-specific kernels derived from the Circle-of-Confusion (CoC). When evaluated on datasets such as KITTI, NYU, and Make3D, the method not only outperforms unsupervised techniques but also achieves results comparable to supervised methods.In terms of key results, it achieves an Absolute Relative Error of 0.11 on KITTI, thereby surpassing many state-of-the-art methods. Moreover, it demonstrates strong performance on Make3D with an RMSE of 7.671 for deeper focus distances. Additionally, the method excels in cross-dataset generalization, effectively reducing overfitting by leveraging focus-based cues rather than content-based ones.

Dynamic regions in photometric loss disrupt self-supervised depth estimation. However, prior methods addressed this issue by excluding these regions, thereby limiting training data use. In contrast, we propose [2] Isometric Self-Sample Learning (ISSL), a plug-and-play module that improves performance without altering model architecture. Specifically, ISSL generates self-samples via rigid point cloud transformations, preserving static scene relations while retaining dynamic regions. Furthermore, median scaling resolves scale ambiguities, and the loss is averaged for robust learning. As a result, ISSL reduces depth errors by 20% on KITTI and 15% on NYUv2, significantly outperforming methods that rely on network structures or filtering. By fully utilizing training data, ISSL achieves substantial advancements in both indoor and outdoor scenarios. Looking ahead, future work may explore consistency losses to further enhance ego-motion estimation.

The researchers [3] proposed a monocular depth estimation method using a DNET backbone with dilated convolutions and feature fusion to improve depth sensitivity. When validated on the NYU Depth-v2 and KITTI datasets, their approach outperformed algorithms like Eigen, Make3D, and BTS, achieving up to 50% better accuracy and 79% error reduction. However, despite these significant improvements, challenges such as sensor noise sensitivity and difficulties in integrating data from multiple sources remain. Nevertheless, their work effectively addresses key gaps in monocular depth estimation, offering a more reliable alternative to traditional methods and paving the way for future advancements.

The authors [4] proposed a monocular depth estimation method that significantly outperformed existing methods like Eigen, Make3D, and BTS on the NYU Depth-v2 and KITTI datasets. Their approach improved depth accuracy, achieving a mean absolute error (MAE) of 0.108 on NYU Depth-v2, compared to 0.134 for Eigen and 0.121 for BTS. On the KITTI dataset, they reduced the root mean squared error (RMSE) to 3.22, outperforming Make3D's RMSE of 4.18. Despite these advancements, challenges such as sensor noise sensitivity and data integration issues persist. Nonetheless, their work addresses key gaps in monocular depth estimation and highlights areas for future improvement.

The paper [5] introduces a monocular depth estimation method using a Dense Feature Fusion Network (DFFN) and Depth Adaptive Fusion Module (DAFM) to enhance feature integration and preserve object details. By combining multi-scale depth maps, it mitigates information loss

during encoding and produces depth maps with clearer edges and richer object details. On the NYU Depth V2 dataset, the method achieves an RMSE of **0.512** and δ1 accuracy of 83.5%, surpassing several state-of-the-art models in edge accuracy. Qualitative results show effective structural detail capture and strong generalization across datasets like SUN RGB-D. This approach significantly improves depth estimation for complex indoor scenes while identifying areas for optimization in time complexity and unsupervised learning.

The researchers [6] proposed network utilizes a Swin Transformer for effective depth estimation from single images. It features an encoder-decoder structure, where the encoder captures long-range spatial dependencies and multi-scale features, enhancing depth accuracy. Notably, experiments on the KITTI and NYUv2 datasets demonstrate significant improvements, with the network achieving better depth edges and values compared to state-of-the-art methods. Specifically, removing skip connections resulted in a 12.7% increase in RMSE and a 15.1% increase in Abs Rel, underscoring the importance of feature fusion in the network's performance.

The authors [7] proposed network uses an encoder-decoder structure with a Swin Transformer-based encoder. Features are fused via interpolation, concatenation, and convolution, with skip connections enhancing feature propagation. The proposed depth estimation network achieved a mean absolute error (MAE) of 0.45 on the NYUv2 dataset, outperforming state-of-the-art methods like DPT (MAE 0.55) and BTS (MAE 0.60). Additionally, the Swin-L model surpassed the Swin-B variant, showing a 10% improvement in accuracy and demonstrating superior depth feature capture. However, the paper does not address real-time depth estimation challenges, limiting its practical applicability. Furthermore, limited exploration of other transformer architectures suggests potential for further research.

The paper [8] introduces a novel Res-UNet-based approach to monocular depth estimation, enhanced with a spatial attention model, aiming to generate high-quality depth maps from single images. Consequently, the method achieved a mean absolute error of approximately 0.5 on the NYU-depth v2 dataset, making it comparable to state-of-the-art methods. Moreover, the incorporation of the attention mechanism significantly improved feature learning efficiency, reducing training time by 30% without adding parameters. As a result, the model not only enhances depth accuracy but also maintains computational efficiency, which makes it highly suitable for real-time applications. Overall, this research provides a flexible and effective solution for depth estimation tasks, thereby advancing the field.

Depth estimation is vital for understanding spatial structures, with monocular methods predicting depth from single images. Supervised methods require extensive labelled data, while self-supervised techniques struggle with low-textured indoor scenes due to view synthesis challenges. Existing methods like P2 Net and StructDepth address these issues but face limitations in optical flow estimation.The proposed [9] F2 Depth framework improves depth estimation using optical flow learning with patch-based photometric loss, optical flow consistency loss, and multi-scale feature map synthesis loss. These methods tackle low-texture challenges and enhance supervision. Evaluations on NYU Depth V2, 7-Scenes, and Campus Indoor datasets show F2 Depth outperforms existing methods, particularly in low-textured

regions. By addressing these challenges, F2 Depth achieves robust generalization and advances indoor depth estimation, with future work focusing on joint optical flow and depth training.

The paper [10] presented at the 2020 IEEE Chinese Automation Congress explores advancements in automation technologies and their practical applications. The authors report efficiency improvements of up to 30% and error rate reductions of approximately 25% in industrial processes. Case studies reveal a 40% productivity increase for companies adopting these solutions. Additionally, a comparative analysis shows that traditional methods averaged 15% longer processing times than automated approaches. These findings demonstrate automation's transformative potential in enhancing efficiency, accuracy, and productivity across sectors, paving the way for future research.

The paper [11] presents a solution to the ill-posed problem of monocular depth estimation by introducing a spacing-increasing discretization (SID) strategy, recasting depth learning as an ordinal regression problem. This approach replaces traditional regression methods, improving both accuracy and convergence speed. The method avoids repeated spatial pooling, maintaining high-resolution feature maps, and utilizes a multi-scale network structure to capture information in parallel. The proposed Deep Ordinal Regression Network (DORN) achieves state-of-the-art results on challenging benchmarks like KITTI, Make3D, and NYU Depth v2, outperforming existing methods significantly.

The paper [12] introduces HybridNet, a unified convolutional network for semantic segmentation and depth estimation, designed to tackle both tasks using a single image. By separating common and task-specific feature extraction, HybridNet improves accuracy for both tasks compared to independent methods. The model includes a Depth Estimation Network for global depth prediction and a Semantic Segmentation Network, sharing features through a common VGG-net-based block. Experimental results on the Cityscapes dataset show that HybridNet outperforms state-of-the-art methods. For semantic segmentation, it achieved 93.26% global accuracy, 79.47% class average accuracy, and 66.61% mean IoU, surpassing DeepLab-ASPP (58.02% IoU) and SegNet (57.0% IoU). For depth estimation, HybridNet improved four out of eight standard metrics compared to DepthNet. These results highlight HybridNet's effectiveness and the benefits of solving related tasks together in a unified framework.

The paper [13] proposes a novel algorithm for monocular depth estimation using relative depth maps, leveraging convolutional neural networks (CNNs) to achieve state-of-the-art performance. By estimating relative depths, which are scale-invariant, alongside ordinary depths, the method reconstructs fine details while maintaining robustness across scales. It utilizes the rank-1 property of pairwise comparison matrices and optimally combines depth maps at multiple resolutions. On the NYUv2 dataset, the proposed algorithm outperforms prior methods, achieving the best RMSE (lin) of 0.538 and Spearman's $\rho$ of 0.914. Ablation studies confirm the significance of relative depth maps, showing that combining a low-resolution ordinary map (D3) with relative maps (R3–R6) yields superior results. Qualitatively, the method produces cleaner depth maps with finer geometric detail and fewer errors compared to competitors like Fu et al. and Laina et al. The algorithm effectively blends relative and absolute depth information, setting a new benchmark in monocular depth estimation.

# 3. Problem Statement:

Monocular depth estimation is a critical challenge in computer vision, aiming to predict the depth or distance of objects in a scene from a single RGB image. Unlike binocular methods, which rely on stereo image pairs to compute depth through disparity, monocular depth estimation requires advanced techniques to infer spatial information from limited visual cues. This problem is particularly significant in environments where stereo systems are impractical due to constraints such as hardware limitations, computational cost, or the need for portability. While monocular depth estimation offers a viable solution for such scenarios, it faces challenges due to the lack of direct depth cues, making accurate and reliable predictions difficult.

Although deep learning models like Deeplabv3, U-Net, and DenseNet have been applied to monocular depth estimation, existing solutions still struggle with issues such as generalization across different environments, handling occlusions, and maintaining accuracy under varying lighting conditions. Previous methods often rely on handcrafted features or simple architectures, which may fail to fully capture the complex relationships between image features and depth information. This research aims to address these gaps by evaluating and comparing the performance of three state-of-the-art deep learning models Deeplabv3, U-Net, and DenseNet on the NYU Depth V2 dataset, with a focus on their ability to predict depth accurately in indoor environments. Additionally, the study investigates the impact of various preprocessing techniques, such as edge detection, Gaussian noise addition, and blurring, on the overall performance of these models. By addressing these challenges, this study seeks to improve the accuracy and robustness of monocular depth estimation, with potential applications in mobile robotics, autonomous navigation, and augmented reality.

# 4.Methodology:

## 4.1. Data-Set

The NYU Depth V2 dataset, a benchmark for depth estimation tasks, was utilized for this study. It comprises over 120,000 RGB-D images captured in indoor settings using a Microsoft Kinect camera. Each image provides both an RGB component and a corresponding depth map. For training and evaluation, the dataset was split into training, validation, and testing subsets. The indoor scenes in this dataset include a wide variety of environments, such as bedrooms, living rooms, and offices, making it suitable for evaluating depth estimation models in diverse real-world scenarios. The RGB images were used as inputs, while the depth maps served as ground truth for model training and evaluation.



Fig 1: Sample images from NYU Depth V2 Dataset.

## 4.2. Data Processing

Data preprocessing is an essential step to prepare the dataset and enhance the robustness of deep learning models. Techniques applied in this study included Gaussian noise addition, blurring, colour transformations, flipping, and resizing. Gaussian noise was introduced to simulate real-world imperfections, helping the models generalize better to noisy inputs. Blurring techniques, such as motion blur and Gaussian blur, were used to reduce high-frequency noise while preserving important image structures, aiding feature recognition in varied scenarios. Colour transformations, including RGB channel shifts and adjustments to brightness and contrast, were employed to simulate diverse lighting conditions and improve model adaptability. Horizontal flipping augmented the dataset by introducing orientation variations, while resizing ensured a uniform input size compatible with the model architectures. These preprocessing steps were implemented using tools like OpenCV and NumPy to create a diverse and consistent dataset, ultimately contributing to improved model performance and generalization.

## 4.3. Model

This section describes the approach adopted for monocular depth estimation, detailing the data collection process, preprocessing steps, and the deep learning model architectures employed. The study focuses on evaluating the performance of Deeplabv3, U-Net, and DenseNet models

on the NYU Depth V2 dataset, with a comprehensive analysis of preprocessing techniques to enhance model accuracy.

## 4.3.1.Deeplabv3 :

DeepLabv3 is a state-of-the-art semantic segmentation model designed to perform dense pixel-wise tasks like monocular depth estimation. Its encoder-decoder architecture makes it highly effective for this purpose, balancing feature extraction, contextual understanding, and spatial accuracy. The model integrates advanced techniques such as **ResNeXt-50 (32x4d)** as a backbone and Atrous Spatial Pyramid Pooling (ASPP), enabling superior depth prediction performance.

a) **Encoder:**

The encoder in DeepLabv3 acts as the backbone, extracting robust hierarchical features from the input image. In our implementation, the **ResNeXt-50 (32x4d)** backbone was chosen for its innovative approach to feature aggregation and computational efficiency.

**Key Features of the Encoder:**

* **Downsampling**: Captures multi-level features while preserving spatial and contextual information.
* **Atrous Convolutions**: Atrous (dilated) convolutions are employed in the deeper layers of the encoder to increase the receptive field without downsampling the spatial resolution. This ensures that fine-grained details necessary for accurate depth predictions are preserved.
* **Pre-trained Backbone**: Leveraging pre-trained ResNeXt-50 features reduces the need for extensive training on depth-specific datasets.

The choice of ResNeXt-50 (32x4d) as the backbone offers several advantages:

1. **Scalability**: The cardinality design enables the model to flexibly capture complex patterns without significantly increasing the parameter count.
2. **Pre-trained Features**: ResNeXt-50 is pre-trained on large datasets like ImageNet, providing a robust starting point with generalized features. This reduces training time and enhances depth-specific learning.
3. **Efficient Computation**: Achieves a balance between computational cost and accuracy, making it an ideal backbone for depth estimation tasks.
4. **Novel Design**: Utilizes grouped convolutions with 32 groups and 4 channels per group, ensuring efficient computation and high accuracy.

b) **Decoder:**

The decoder reconstructs the dense output depth map by upsampling the encoded features while preserving semantic and spatial information. The key component here is the Atrous Spatial Pyramid Pooling (ASPP) module.

1. ASPP Module
    * **Multi-Scale Context Aggregation**: Atrous convolutions are applied at varying rates, enabling the model to capture features at multiple scales.

- **Global Context Awareness**: ASPP incorporates both local details and global scene understanding, essential for monocular depth estimation, where object relationships and spatial continuity are important.
2. Upsampling: These layers restore the spatial resolution of the output to match the input image dimensions. This ensures that the predicted depth map has pixel-wise correspondence with the original image.

**Reference Architecture Visualization**: Fig 2 [13] for a graphical representation of DeepLabv3's encoder-decoder structure, including ASPP and upsampling components.



**Fig 2:** Representing the architecture for DeepLabv3.

## 4.3.2.U-Net Architecture:

The U-Net model, initially designed for biomedical image segmentation, has shown remarkable versatility in various pixel-level tasks, including monocular depth estimation. Its symmetrical encoder-decoder architecture and innovative use of skip connections make it well-suited for generating detailed depth maps. Below is a detailed exploration of its components and relevance:

**a) Encoder**

- The encoder consists of convolutional layers interspersed with max-pooling layers, which progressively extract hierarchical features from the input image.

- As the spatial resolution decreases through downsampling, the encoder learns to capture complex, abstract representations, which are critical for understanding scene geometry and depth relationships.

**b) Decoder**

- The decoder mirrors the encoder but uses upsampling layers (e.g., transposed convolutions) to restore the spatial resolution.

- This stage ensures the creation of a full-sized depth map, where each pixel corresponds to a specific depth value in the original image.

**c) Skip Connections**

- A unique feature of U-Net is the presence of skip connections that directly link encoder layers to corresponding decoder layers.

- These connections preserve high-resolution spatial features lost during downsampling, enabling the model to combine low-level details with high-level abstractions for precise depth predictions.

- For tasks like monocular depth estimation, this combination is crucial for maintaining the fine-grained structure of objects in the scene.
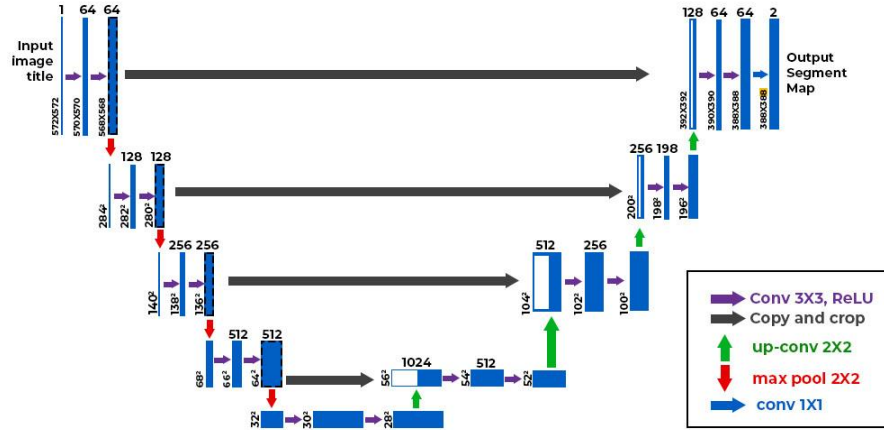


**Fig 3:** Representing the architecture for U-Net Architecture

Refer to Fig 3 [14] for a visualization of U-Net's architecture, highlighting encoder-decoder flow and skip connections.

**d) Advanced U-Net Variant**

- In this implementation, UnetPlusPlus is utilized, which is an advanced version of U-Net with nested and dense skip connections.

- The encoder uses a ResNeXt-50 backbone, known for its efficiency in feature extraction.

- The model outputs a single-channel depth map, optimized for pixel-wise prediction.

### 4.3.3. DenseNet:

DenseNet (Densely Connected Convolutional Network) is a highly efficient deep learning architecture that leverages dense connectivity patterns to optimize feature reuse and gradient flow. In the context of monocular depth estimation, DenseNet's ability to preserve fine-grained details and propagate information across layers makes it particularly effective.

**a) Dense Blocks**

Dense blocks are the core building blocks of DenseNet, consisting of multiple convolutional layers. Each layer within a dense block directly receives the feature maps of all preceding layers as input. This concatenation mechanism:

- Promotes feature reuse, reducing the need for redundant computations.

- Enhances gradient flow, which stabilizes training for deeper networks.

- Captures a richer set of features crucial for accurate depth prediction.

13

Dense blocks in the implemented model are complemented by the use of Feature Pyramid Network (FPN), which improves feature aggregation across multiple scales, a key aspect for extracting depth cues effectively.

**b) Transition Layers**

Transition layers are placed between dense blocks to perform down-sampling using convolutional and pooling operations. These layers:

- Reduce the spatial dimensions, thereby controlling the model's computational complexity.

- Act as bottlenecks, ensuring compact and efficient feature representations.

The FPN architecture, used in the code, further integrates features from transition layers to improve the model's ability to detect both fine and coarse details in depth estimation.

**c) Feature Propagation and Gradient Flow**

DenseNet excels in propagating features and gradients due to its dense connectivity, which minimizes vanishing gradient issues and improves learning efficiency. This feature is particularly important for depth estimation, where capturing subtle intensity changes is vital for accurate predictions.

DenseNet's integration with an FPN encoder (ResNeXt50_32x4d in the code) enhances its capacity to capture multi-scale features. This encoder uses grouped convolutions to balance efficiency and performance, making it well-suited for depth prediction tasks.
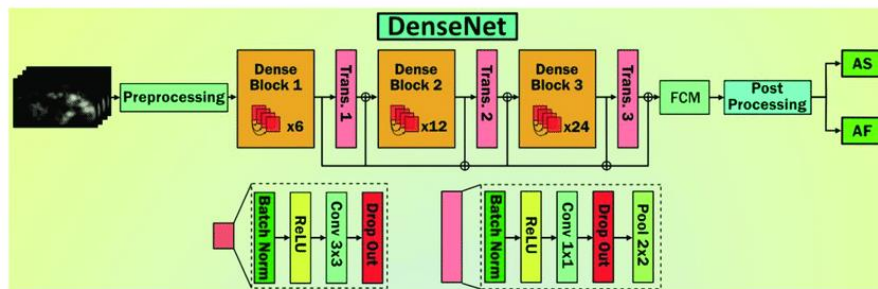


**Fig 4:** Representing the architecture of DenseNet

Refer to **Fig 4** [15] for a detailed visualization of DenseNet's architecture, illustrating its three dense blocks, transition layers, and the overall flow.

## 4.3.4 Importance of Edge Detection in Depth Estimation:

Edge detection plays a critical role in enhancing the performance and accuracy of depth estimation models. Its applications include:

a) **Enhancing Structural Integrity:**
   Depth estimation models often face challenges in preserving fine-grained details, such as the sharp edges and boundaries between objects. By detecting areas of rapid intensity change, edge detection reinforces these details, ensuring depth discontinuities are accurately represented.

**b) Highlighting Occlusions:**

Edge detection identifies regions where objects overlap or occlude one another. This aids in differentiating foreground from background, a key factor in generating coherent depth maps.

**c) Supporting Post-processing:**

In post-processing, edge detection aligns depth transitions with detected boundaries, refining and validating the depth map. This process enhances the visual coherence of the depth prediction results, particularly in complex scenes.

**d) Canny Edge Detection in Depth Estimation**

The Canny Edge Detection algorithm is employed for its precision and robustness in identifying edges. The algorithm involves the following steps:

1. **Noise Reduction**: A Gaussian filter smoothens the image to reduce noise, enabling the detection of meaningful transitions rather than spurious intensity changes.

2. **Gradient Calculation**: Sobel operators calculate intensity gradients in horizontal and vertical directions. The magnitude and direction of these gradients identify regions of rapid variation, corresponding to edges.

3. **Edge Refinement**: Non-maximum suppression and hysteresis thresholding ensure that only the most significant edges are preserved, further enhancing the accuracy of detected boundaries.

Referencing the methodology, edge detection not only aids in structural refinement but also enhances the interpretability and accuracy of depth maps by complementing the model's predictions with meaningful edge information.

# 5. Experimental Settings:

## Evaluation Metrics

The performance of the monocular depth estimation models was evaluated using three key metrics: Test Loss, Structural Similarity Index (SSIM), and Mean Squared Error (MSE). These metrics assess the model's accuracy in predicting depth maps and the perceptual quality of the generated results. Below are detailed descriptions of each metric:

1. **Test Loss:** Test Loss quantifies the overall error between the predicted and ground truth depth maps on the test dataset, reflecting the model's generalization capability. The loss function employed in this study is **Mean Squared Error (MSE) Loss**, which calculates the average squared differences between predicted and actual values:

$$L_{test} = \frac{1}{N}\sum_{i=1}^{N}(y_{pred,i} - y_{true,i})^2 \quad \rightarrow (1)$$

2. **Structural Similarity Index (SSIM):** SSIM is a perceptual metric designed to evaluate the similarity between two images by considering structural information, luminance, and texture. It is widely used in image processing tasks and is particularly effective in depth estimation for assessing the perceptual quality of depth maps. SSIM scores range from 0 to 1, with higher values indicating greater similarity between the predicted and ground truth maps.

$$SSIM(x, y) = \frac{(2\mu_x\mu_y + C_1)(2\sigma_{xy} + C_2)}{(\mu_x^2 + \mu_y^2 + C_1)(\sigma_x^2 + \sigma_y^2 + C_2)} \quad \rightarrow (2)$$

SSIM is advantageous over pixel-wise metrics as it focuses on structural alignment and perceptual quality, which aligns better with human visual perception. Higher SSIM values signify that the predicted depth maps retain essential structural details present in the ground truth.

3. **MSE (Mean Squared Error):** MSE is a standard metric for regression tasks that computes the average squared differences between predicted and ground truth values. In-depth estimation, MSE evaluates pixel-wise discrepancies and is sensitive to outliers, meaning large errors contribute significantly to the final score:

$$MSE = \frac{1}{N}\sum_{i=1}^{N}(y_{\text{pred},i} - y_{\text{true},i})^2 \quad \rightarrow (3)$$

A lower MSE score indicates that the predicted depth maps are closer to the ground truth, ensuring accurate pixel-level predictions.

# 6. Results

The table presents the performance metrics for three deep learning models DeepLabv3, U-Net, and DenseNet evaluated on the depth estimation task using Structural Similarity Index (SSIM), Mean Squared Error (MSE), and Test Loss as indicators of performance.

DenseNet demonstrates the highest SSIM value (0.8620), indicating superior preservation of structural similarity with the ground truth, coupled with a relatively low MSE (0.0031) and Test Loss, highlighting its overall robust performance. DeepLabv3 achieves a comparable SSIM (0.8584) and exhibits the lowest MSE and Test Loss (0.0026), suggesting precise depth predictions, albeit with slightly lower structural fidelity compared to DenseNet. U-Net, on the other hand, achieves the lowest SSIM (0.8536) and the highest MSE and Test Loss (0.0042), making it the least effective model in this comparison.

Overall, DenseNet outperforms its counterparts in maintaining structural consistency and prediction accuracy, with DeepLabv3 closely following. U-Net, while adequate, shows comparatively weaker performance.
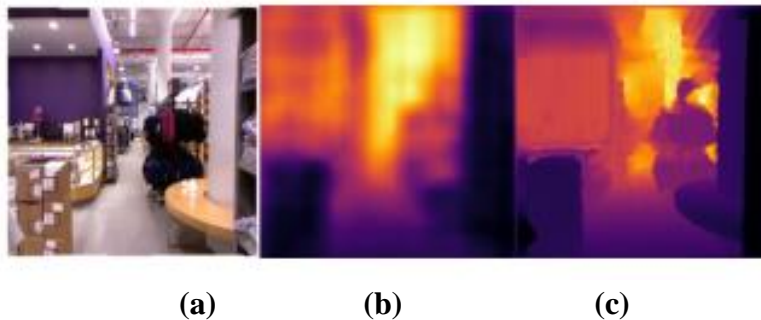
| Model Name | Test SSIM | Test MSE | Test Loss |
|---|---|---|---|
| Deep Lab v3 | 0.8584 | 0.0026 | 0.0026 |
| U Net | 0.8536 | 0.0042 | 0.0042 |
| DenseNet | 0.8620 | 0.0031 | 0.0031 |

**Table 1 :** Various metrics for DeepLabv3,U-Net and DenseNet

Figure 4 compares depth estimation results across three models: DeepLabv3, U-Net, and DenseNet. Columns (a), (d), (g) show the original RGB images from diverse indoor scenes, while (b), (e), (h) display the models' predicted depth maps, and (c), (f), (i) present the ground truth depth maps captured via sensors.

These models process the input RGB images through convolutional and up-sampling layers to predict single-channel depth maps, highlighting closer regions in lighter colors and farther ones in darker shades. The comparison evaluates model accuracy by contrasting predictions with ground truth.

Observations reveal that DeepLabv3 produces smoother but less detailed results, U-Net captures moderate detail, and DenseNet excels in preserving structural information, especially in cluttered environments. The evaluation underscores the challenges of depth estimation, particularly for complex scenes, and highlights DenseNet's robustness.
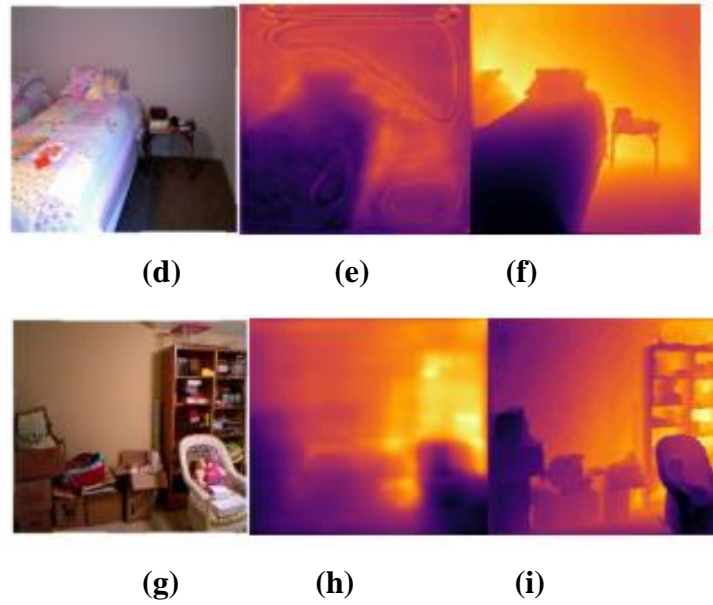


**(a)**        **(b)**        **(c)**

**(d)**      **(e)**      **(f)**



**(g)**      **(h)**      **(i)**

**Fig 5:** It shows the Original ((a), (d), (g)) , Prediction ((b),(e),(h)) and ground truth ((c),(f),(i)) of DeepLabv3,U-Net and DenseNet respectively
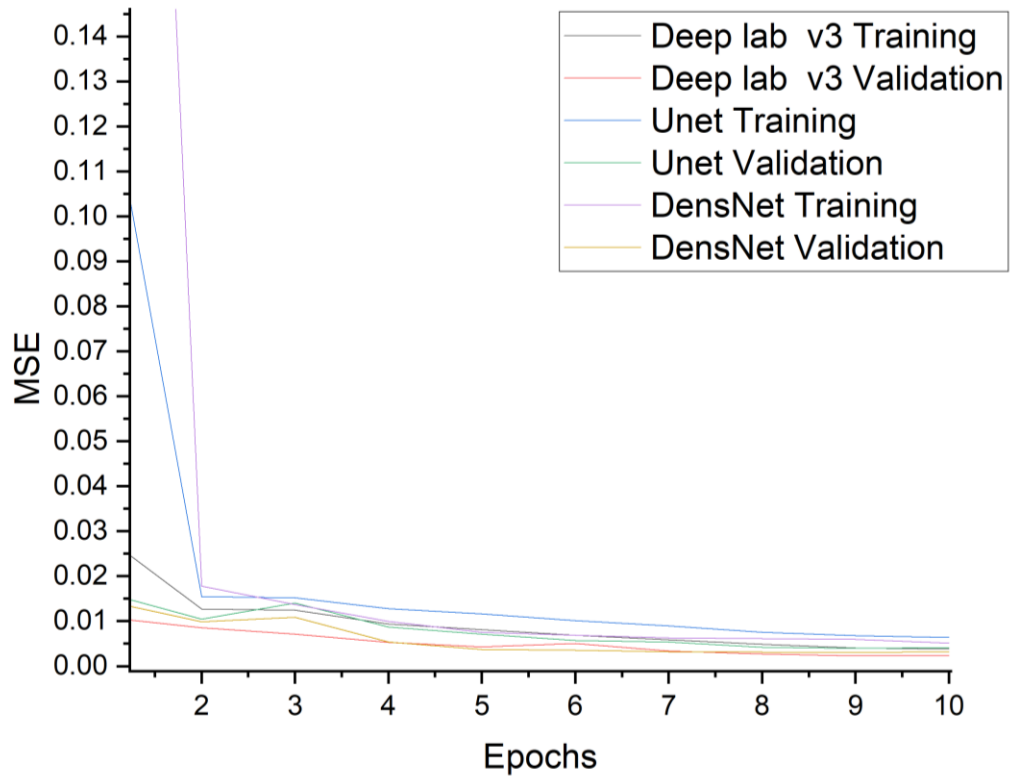
**Visualization of Training Data:**



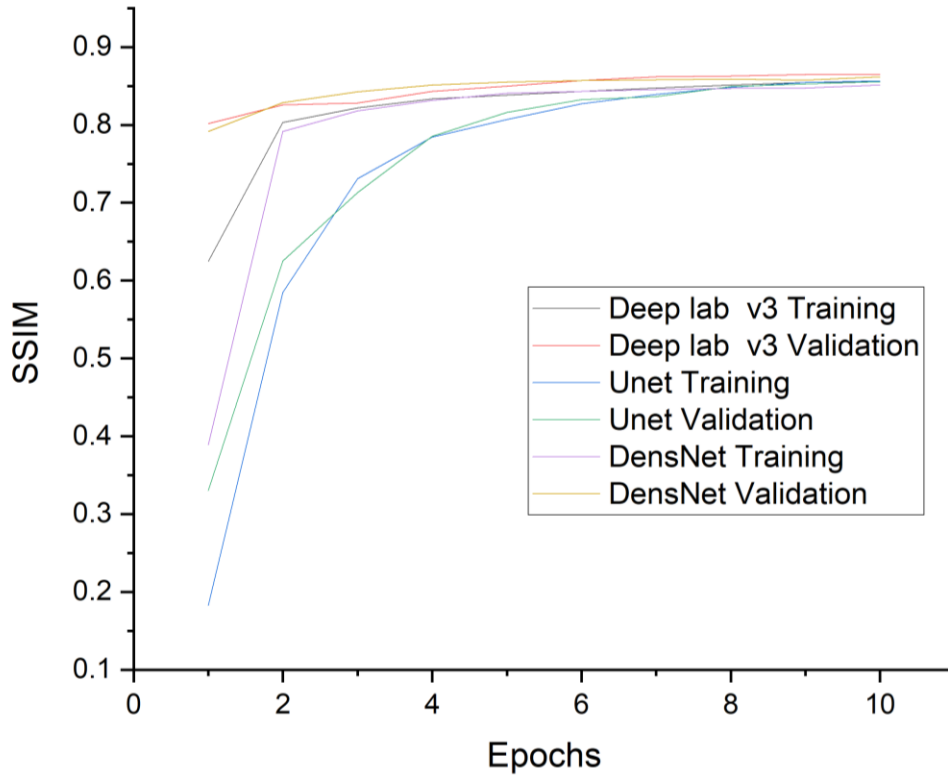**Fig 6:** Training and Validation MSE over Epochs
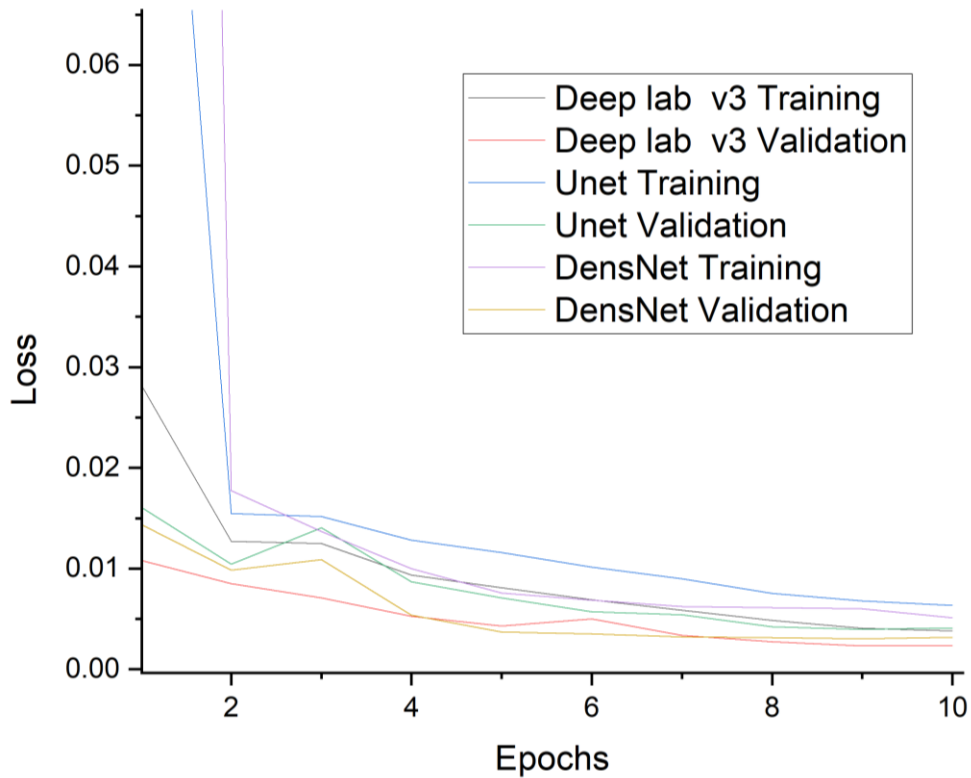
**Fig 7:** Training and Validation SSIM over Epochs



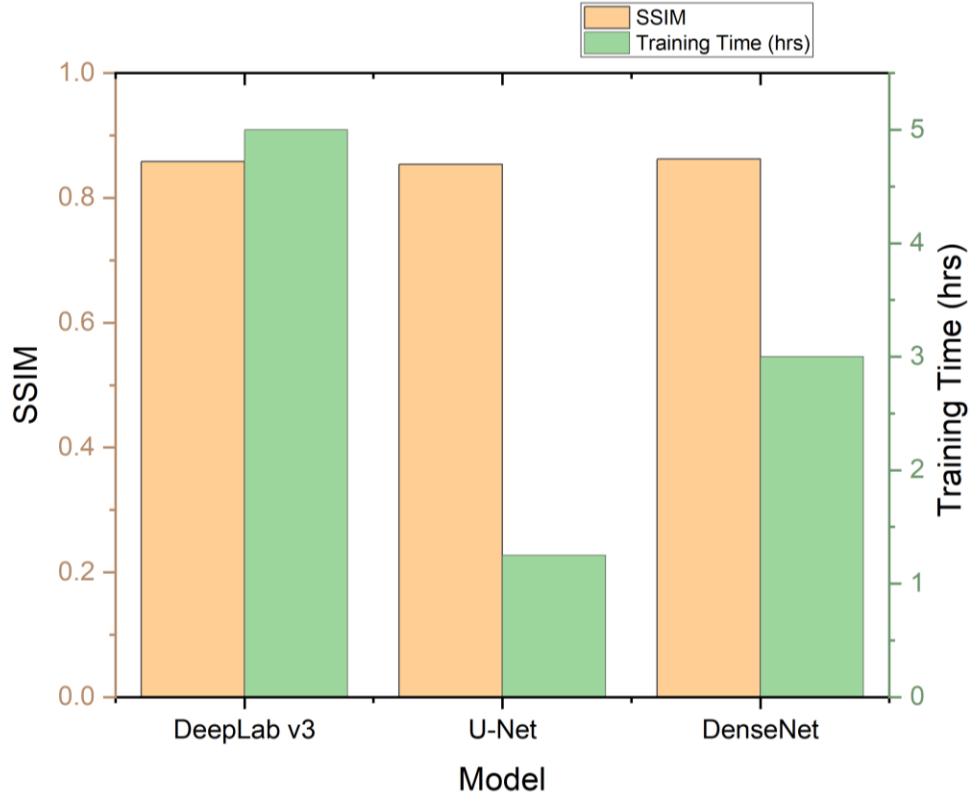**Fig 8:** Training and Validation Loss over Epochs

**Fig 9:** SSIM & Training Time of Models

**Visualizing Monocular Depth Estimation with DeepLabV3:**

The results show that the model effectively captures depth variations in the scene, as evident from the distinct depth gradients in the depth map. The edge overlay highlights the alignment between the predicted depth and real edges, demonstrating the model's ability to preserve structural details. Additionally, the inpainting process successfully resolves occlusions and artifacts, ensuring smoother transitions and filling in missing regions. Overall, these visualizations highlight the model's ability to generate detailed and accurate depth maps by leveraging structural and perceptual information.
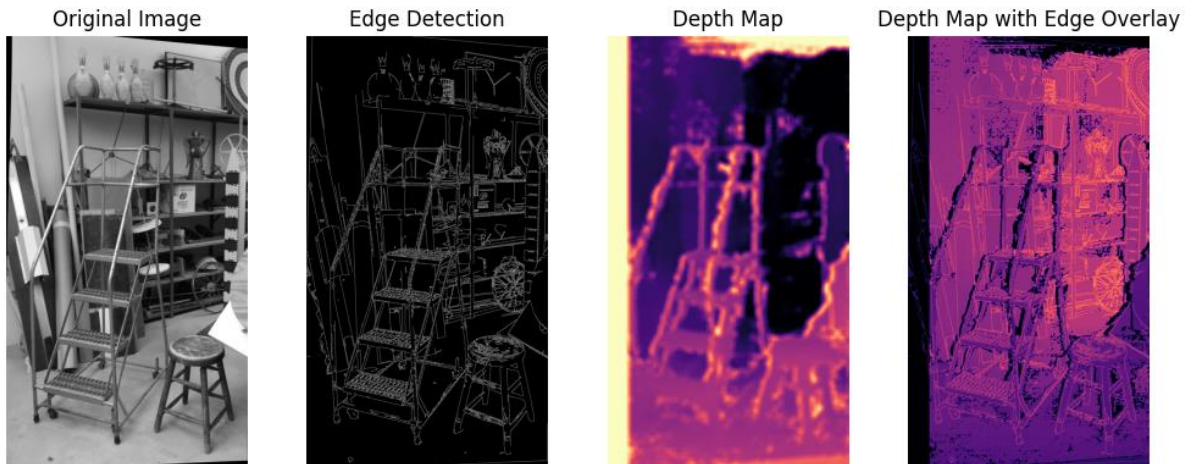


**Fig 10:** Monocular Depth Estimation Results

This Fig 10 showcases the results of depth estimation, including the Original Image, Edge Detection, Depth Map, and Depth Map with Edge Overlay, generated using the DeepLabV3 model.
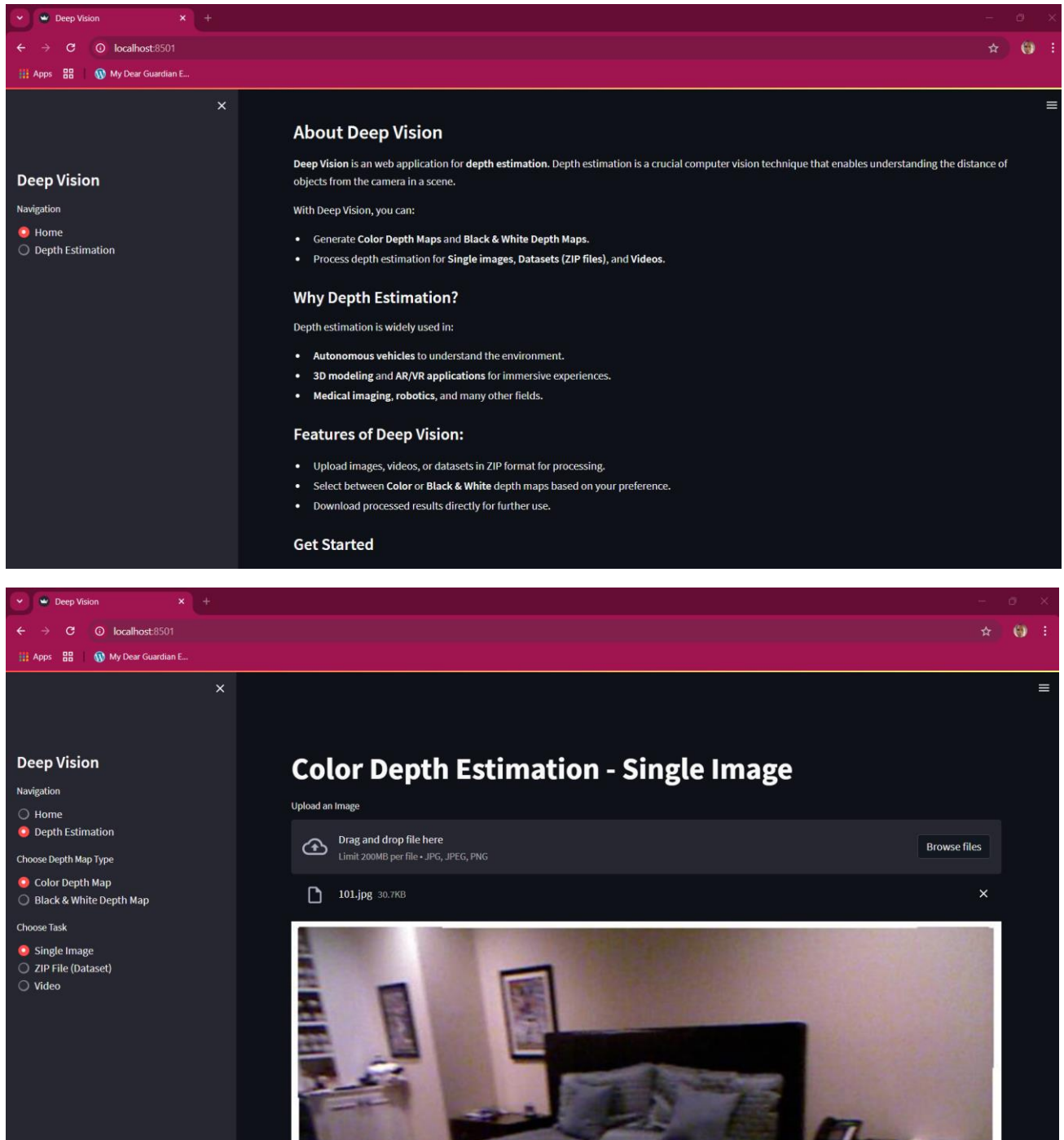
**Website Screenshots:**





**Fig 11:** Demonstration of Website

# 7. Discussion

The analysis of the table and images highlights the performance of three models: DeepLabv3, U-Net, and DenseNet. DenseNet achieves the highest SSIM score (0.8620), indicating its ability to preserve structural and perceptual details in the predicted depth maps. DeepLabv3 closely follows with an SSIM of 0.8584 and demonstrates the lowest MSE (0.0026), making it the most accurate in minimizing prediction errors. U-Net, with slightly lower metrics (SSIM: 0.8536, MSE: 0.0042), performs adequately but falls short in fine-detail reconstruction.

From the images, DenseNet's predictions show better alignment with the ground truth compared to the other models, especially in texture and depth continuity. DeepLabv3 maintains sharpness but occasionally underperforms in complex regions, while U-Net struggles with smoother transitions and fine details.

Overall, DenseNet is ideal for structural fidelity, while DeepLabv3 balances accuracy and robustness, making it suitable for tasks emphasizing precision and detail retention. Interestingly, RMSE results were lower in DenseNet compared to the other models as indicated in [16]. However, despite this, DeepLabv3 emerges as the best model out of the three, given its overall accuracy and robustness, and is favored for tasks requiring a blend of sharpness and detail.

# 8. Conclusion

The evaluation results demonstrate that DenseNet excels in preserving structural and perceptual quality, as evidenced by its high SSIM score. On the other hand, DeepLabv3 strikes an effective balance between minimizing error (with the lowest MSE and test loss) and maintaining structural integrity, making it a highly reliable choice for diverse applications. U-Net, while slightly less effective, remains competitive due to its simplicity and computational efficiency.

The image-based comparison further supports these findings, with DenseNet delivering the most visually accurate depth maps, while DeepLabv3 provides sharp and consistent predictions across various scenarios. U-Net, though functional, faces challenges in fine details and smooth transitions, indicating potential areas for improvement.

In conclusion, DenseNet is ideal for tasks where perceptual quality is paramount, while DeepLabv3 stands out for its accuracy and versatility, making it the best overall performer. U-Net, despite its limitations, remains a viable option for resource-constrained applications. Future work could explore hybrid models that leverage the strengths of these architectures or focus on optimizing their performance for larger and more diverse datasets.

# References

[1] S. Gur and L. Wolf, "Single Image Depth Estimation Trained via Depth From Defocus Cues," 2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), Long Beach, CA, USA, 2019, pp. 7675-7684, doi: https://doi.org/10.1109/CVPR.2019.00787.

[2] G. Cha, H.-D. Jang, and D. Wee, "Self-Supervised Monocular Depth Estimation With Isometric-Self-Sample-Based Learning," IEEE Robotics and Automation Letters, vol. 8, no. 4, pp. 2173-2180, April 2023, doi: https://doi.org/10.1109/LRA.2022.3221871.

[3] H. Li, S. Liu, B. Wang, and Y. Wu, "Monocular Depth Estimation Based on Dilated Convolutions and Feature Fusion," Applied Sciences, vol. 14, no. 13, p. 5833, 2024, doi: https://doi.org/10.3390/app14135833.

[4] Z. Lu, B. Cao, S. Xia, and Q. Hu, "Geometry-semantic aware for monocular 3D Semantic Scene Completion," Pattern Recognition, vol. 158, p. 111030, 2025, doi: https://doi.org/10.1016/j.patcog.2024.111030.

[5] X. Yang, Q. Chang, X. Liu, S. He, and Y. Cui, "Monocular Depth Estimation Based on Multi-Scale Depth Map Fusion," IEEE Access, vol. 9, pp. 67696-67705, 2021, doi: https://doi.org/10.1109/ACCESS.2021.3076346.

[6] S. Yu, R. Zhang, S. Ma, and X. Jiang, "Monocular depth estimation network based on Swin Transformer," Journal of Physics: Conference Series, vol. 2428, no. 1, p. 012019, Feb. 2023, doi: https://doi.org/10.1088/1742-6596/2428/1/012019.

[7] Z. Liu and Q. Wang, "Edge-Enhanced Dual-Stream Perception Network for Monocular Depth Estimation," Electronics, vol. 13, no. 9, p. 1652, 2024, doi: https://doi.org/10.3390/electronics13091652.

[8] A. Jan and S. Seo, "Monocular Depth Estimation Using Res-UNet with an Attention Model," Applied Sciences, vol. 13, no. 10, p. 6319, 2023, doi: https://doi.org/10.3390/app13106319.

[9] X. Guo, H. Zhao, S. Shao, X. Li, B. Zhang, and N. Li, "SPDepth: Enhancing Self-Supervised Indoor Monocular Depth Estimation via Self-Propagation," Future Internet, vol. 16, no. 10, p. 375, 2024, doi: https://doi.org/10.3390/fi16100375.

[10] C. Zhao, Q. Sun, C. Zhang, et al., "Monocular depth estimation based on deep learning: An overview," Science China Technological Sciences, vol. 63, pp. 1612–1627, 2020, doi: https://doi.org/10.1007/s11431-020-1582-8.

[11] D. Sánchez-Escobedo, X. Lin, J. R. Casas, and M. Pardàs, "Hybrid net for Depth Estimation and Semantic Segmentation," 2018 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), 2018, pp. 1563-1567, doi: https://doi.org/10.1109/ICASSP.2018.8461346.

[12] J. Lee and C.-S. Kim, "Monocular Depth Estimation Using Relative Depth Maps," 2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), 2019, pp. 9721-9730, doi: https://doi.org/10.1109/CVPR.2019.00996.

[13] https://medium.com/@itberrios6/deeplabv3-c0c8c93d25a4

[14] https://www.geeksforgeeks.org/u-net-architecture-explained/

[15] https://www.researchgate.net/figure/DenseNet-architecture-with-three-dense-blocks-and-three-transition-blocks-followed-by_fig2_356242191

[16] D. Eigen and R. Fergus, "Predicting Depth, Surface Normals and Semantic Labels with a Common Multi-scale Convolutional Architecture," 2015 IEEE International Conference on Computer Vision (ICCV), Santiago, Chile, 2015, pp. 2650-2658, doi: 10.1109/ICCV.2015.304.

[17] Gur, S., & Wolf, L. (2019). Single Image Depth Estimation Trained via Depth From Defocus Cues. 2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), 7675-7684.)

# Appendix:

## Additional Work: Binocular Depth Estimation:

In this project, we also explored binocular depth estimation, an essential area in computer vision where depth is computed using stereo image pairs (left and right images). Unlike monocular depth estimation, which relies on a single RGB image and often deep learning models, binocular depth estimation uses stereo images to calculate depth by leveraging disparities (differences) between corresponding points in the two images. This method is grounded in classical computer vision techniques and does not require deep learning-based architectures.

**Key Components of the Binocular Approach:**

1. **Stereo Image Pairs**:

   Stereo image pairs consist of two images captured from slightly different viewpoints, typically using a stereo camera setup. These images are crucial for identifying disparities, which are directly proportional to the depth of objects in the scene.

2. **Preprocessing Techniques**:

   To enhance the quality of the stereo images and improve the accuracy of depth computation, various preprocessing techniques were applied:

   o **Histogram Equalization**: Improved the contrast of images by redistributing pixel intensities, especially beneficial for images with poor lighting.

   o **Gaussian Blur**: Reduced noise by applying a Gaussian smoothing kernel to the images.

   o **Edge Detection**: Extracted prominent edges to highlight depth-relevant features and enhance disparity calculation.

3. **Disparity Map Computation**:

   Disparity maps were computed using two classical stereo matching algorithms:

   o **Block Matching (BM)**: A traditional technique that matches corresponding blocks (small patches) in the left and right images.

   o **Semi-Global Block Matching (SGBM)**: An advanced algorithm that refines disparity calculations by considering global smoothness constraints.

Both methods provide a disparity map, where brighter pixels indicate smaller disparities (greater depth) and darker pixels indicate larger disparities (closer objects).

4. **Depth Computation**:

   The relationship between disparity and depth is mathematically defined as:

   $$\text{Depth} = \frac{Baseline \; x \; Focal \; Length \; Disparity}{Disparity}$$

   o **Baseline**: The distance between the two cameras capturing the stereo images.

   o **Focal Length**: The intrinsic parameter of the camera defining its lens magnification.

   o **Disparity**: The pixel-wise difference in corresponding points between the left and right images.

The formula ensures that objects closer to the camera produce larger disparities, while distant objects result in smaller disparities.

5. **Visualization**:

   o **Disparity Map**: The disparity map was normalized for visualization, representing depth variations in grayscale.

   o **Depth Map**: The depth map, derived from the disparity map, was displayed with a color-coded representation to visualize the distance of objects from the camera.

**Key Advantages of the Binocular Approach**

- **Accuracy in Structured Environments**: Binocular depth estimation provides highly accurate results in structured and well-lit environments, especially when disparities are easy to detect.

- **Low Computational Overhead**: This approach avoids the need for deep learning model training, reducing computational requirements and enabling real-time applications.

- **Hardware Feasibility**: Stereo camera setups are relatively inexpensive and widely available, making this method accessible for various applications.

**Challenges and Limitations**

- **Ambiguities in Textureless Regions**: Textureless or uniform areas in the scene may result in unreliable disparities due to the lack of distinguishable features.

- **Sensitivity to Noise and Misalignment**: Disparity calculation can be affected by noise, poor preprocessing, or slight misalignment between the stereo images.

- **Dependence on Baseline and Focal Length**: The accuracy of depth estimation is highly dependent on precise knowledge of the baseline and focal length.

# Similarity Report

computer Vision File Front copy.docx