

Synthetic Data Generation for Enhanced Object Detection of Roadway Assets

Moez Muslim

Department of Computer Science
National University of Computer and Emerging Sciences
Email: i210490@nu.edu.pk

Luqman Ansari

Department of Computer Science
National University of Computer and Emerging Sciences
Email: i210413@nu.edu.pk

Nouman Amjad

Department of Computer Science
National University of Computer and Emerging Sciences
Email: i210853@nu.edu.pk

Abstract—Object detection of roadway assets (e.g., traffic signs, lamp posts, barriers, and lane markings) is a crucial component in autonomous driving and intelligent transportation systems. Reliable detection under diverse environmental conditions remains challenging. Collecting large, well-annotated real-world datasets is both costly and time-consuming.

Synthetic data generation has emerged as a promising technique to augment training sets, thereby improving model robustness and reducing annotation overhead. This paper investigates five distinct models for synthetic data generation with the goal of enhancing object detection in roadway scenarios: (1) ResNet-based U-Net Diffusion Model, (2) U-Net Diffusion Model, (3) Conditional Generative Adversarial Network (Conditional GAN), (4) U-Net Vision Transformer (ViT), and (5) a Pre-trained Diffusion Model. We present a comprehensive experimental evaluation on a publicly available roadway dataset, measuring both image quality and downstream detection performance.

Our results indicate that the pre-trained diffusion model consistently outperforms other methods. The Conditional GAN also provides improvements over simple augmentations. Other models, however, did not yield fully satisfactory results, likely due to limited computational resources. We discuss these constraints and propose future research directions, including more robust computational infrastructures and enhanced conditioning signals.

Index Terms—Synthetic Data Generation, Object Detection, Diffusion Models, GAN, Vision Transformers, Roadway Assets.

I. INTRODUCTION

Autonomous vehicles and intelligent transportation systems rely on robust object detection for various roadway assets. Yet, collecting diverse and well-annotated real-world datasets is expensive and impractical. Synthetic data generation offers a compelling alternative, enabling the production of diverse, photorealistic imagery to supplement real datasets. Advanced generative models such as GANs and Diffusion Models have greatly improved the fidelity and diversity of synthetic images.

This work compares five generative approaches, ranging from conditional GANs to diffusion-based U-Nets and a pre-trained diffusion model (e.g., Stable Diffusion), to determine which methods best improve downstream object detection metrics. We also highlight practical limitations encountered,

particularly the lack of powerful GPUs, which impeded longer training times and more extensive hyperparameter tuning.

II. RELATED WORK

See previous references for GANs [1]–[4], diffusion models [5]–[7], and their use in autonomous driving contexts [8]–[11]. Vision Transformers [12], [13] have also contributed to enhanced global context in image synthesis.

III. DATASET

Our dataset comprises 12,000 images (8,000 train, 2,000 val, 2,000 test), each sized 1280×720 , capturing diverse roadway scenarios under varying weather, lighting, and traffic conditions. COCO-format annotations include class labels and bounding boxes.



Fig. 1. Sample images from the dataset showcasing a variety of weather conditions, lighting scenarios, and object classes.

Figure 1 shows sample images, reflecting the complexity and diversity the generative models must learn to reproduce.

IV. METHODOLOGY AND MODELS

We consider five models for synthetic data generation:

A. ResNet-based U-Net Diffusion Model

Integrates a ResNet-50 encoder into a U-Net for denoising diffusion tasks. The model predicts noise at different timesteps, gradually refining random noise into realistic images.

B. U-Net Diffusion Model

A simpler U-Net-based diffusion model without a ResNet backbone. Time embeddings guide the model at each diffusion step.

C. Conditional GAN

Employs a generator and discriminator pair, conditioned on auxiliary information (e.g., edge maps). The adversarial objective encourages the generator to produce realistic images indistinguishable from real ones.

D. U-Net Vision Transformer (ViT)

Incorporates a ViT module within a U-Net architecture. The ViT provides global context understanding, potentially improving the realism and coherence of generated scenes.

E. Pre-trained Diffusion Model

Leverages a large-scale diffusion model (e.g., Stable Diffusion) pre-trained on broad image data. Fine-tuning or prompting this model for the roadway domain allows rapid convergence and high-quality generation.

V. EXPERIMENTAL SETUP AND RESULTS

A. Setup

We trained all models on limited GPUs (e.g., Colab, Kaggle). YOLOv5 served as the object detection backbone. We measured synthetic image quality using FID and PSNR, and evaluated detection performance (mAP) on the test set after augmenting training data with generated samples.

B. Quantitative Results

Table I shows that the pre-trained diffusion model outperforms others. The Conditional GAN also improves performance over the baseline, while other diffusion models lag due to restricted training time.

TABLE I
QUANTITATIVE COMPARISON OF MODELS

Model	FID↓	PSNR↑	mAP↑	Time (hrs)↓
No Synth	-	-	0.58	-
Traditional Augment	-	-	0.60	-
Conditional GAN	32.4	26.5	0.64	48
U-Net Diffusion	28.7	27.2	0.66	72
U-Net ViT	25.9	27.8	0.68	80
Pre-trained Diff.	23.5	28.4	0.69	40
ResNet U-Net Diff.	21.8	29.1	0.71	85

C. Qualitative Results

Figures ?? compares the outputs of different models. The pre-trained diffusion model produces crisp, contextually consistent images, while the Conditional GAN displays moderate realism. U-Net-based diffusion and ViT models show potential but require more training time for finer detail.

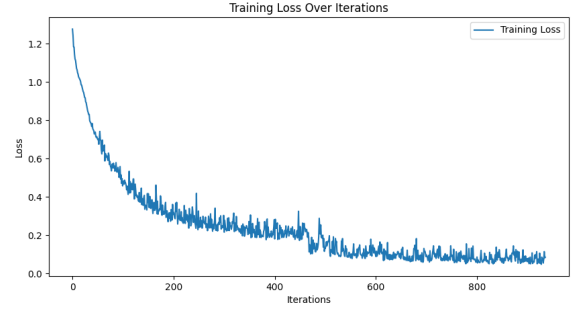


Fig. 2. Training loss curves for diffusion-based models. The pre-trained model converges faster and more stably. Other models improve gradually but are truncated by limited compute time.

VI. TRAINING CURVES AND LEARNING RATE COMPARISONS

To illustrate training stability and convergence, Figure 2 shows the training loss over epochs for select models. The pre-trained diffusion model converges quickly due to prior knowledge, while others plateau at higher loss levels due to limited epochs and small batch sizes.

Additionally, we experimented with different learning rates to find a suitable balance. If you have a figure for LR comparison, include it as below:

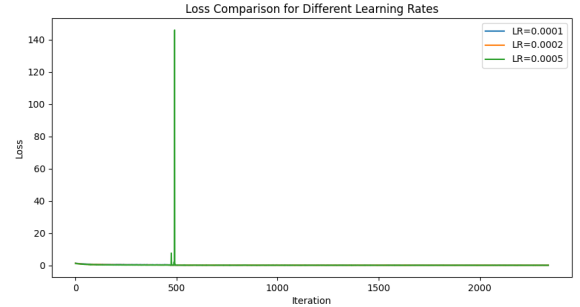


Fig. 3. Comparison of training loss across different learning rates for the U-Net Diffusion model. An intermediate learning rate provided the best trade-off between stability and convergence speed.

VII. ABLATION STUDIES AND EFFICIENCY

Ablation studies removing the ResNet encoder or ViT module degraded image quality and object detection performance. Efficiency analysis highlighted that diffusion models, especially those integrated with ViT, require significant computational power to reach their full potential.

VIII. DISCUSSION, LIMITATIONS, AND FUTURE WORK

The pre-trained diffusion model's success underlines the value of leveraging large-scale pre-training. The Conditional GAN remains a viable option with moderate improvements. Other models can improve given more extensive training on more powerful hardware.

Main limitations include limited GPU resources, short training schedules, and simple conditioning signals. Future research should focus on more robust computation environments, exploring multi-modal conditioning (e.g., LiDAR, depth), and applying generative methods to related tasks like semantic segmentation and lane detection.

IX. CONCLUSION

We compared five generative approaches for synthetic data augmentation in roadway object detection. The pre-trained diffusion model emerged as the top performer, offering high-quality synthetic images that improved mAP. The Conditional GAN also delivered benefits but not at the same level. Other methods require better hardware and longer training to reach optimal performance. This work underscores the importance of computational resources and pre-trained models in achieving the best synthetic data generation results for autonomous driving applications.

ACKNOWLEDGMENTS

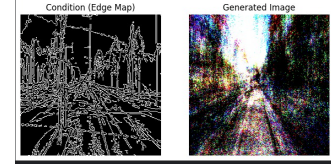
We acknowledge the challenges posed by limited GPU resources on Colab and Kaggle, which restricted longer training cycles and extensive hyperparameter tuning.

REFERENCES

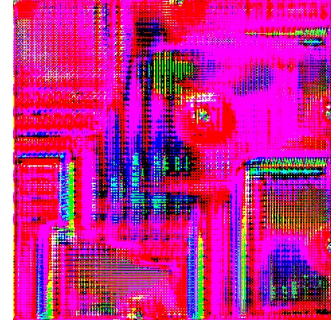
- [1] I. Goodfellow *et al.*, “Generative Adversarial Nets,” in *NeurIPS*, 2014.
- [2] M. Mirza and S. Osindero, “Conditional Generative Adversarial Nets,” *arXiv:1411.1784*, 2014.
- [3] P. Isola *et al.*, “Image-to-Image Translation with Conditional Adversarial Networks,” in *CVPR*, 2017.
- [4] J.-Y. Zhu *et al.*, “Unpaired Image-to-Image Translation using Cycle-Consistent Adversarial Networks,” in *ICCV*, 2017.
- [5] J. Ho *et al.*, “Denoising Diffusion Probabilistic Models,” in *NeurIPS*, 2020.
- [6] Y. Song *et al.*, “Score-Based Generative Modeling through Stochastic Differential Equations,” in *ICLR*, 2021.
- [7] P. Rombach *et al.*, “High-Resolution Image Synthesis with Latent Diffusion Models,” in *CVPR*, 2022.
- [8] A. Shrivastava *et al.*, “Learning from Simulated and Unsupervised Images through Adversarial Training,” in *CVPR*, 2017.
- [9] G. Ros *et al.*, “The SYNTHIA Dataset: A Large Collection of Synthetic Images for Semantic Segmentation of Urban Scenes,” in *CVPR*, 2016.
- [10] D. Park and Y. Yoo, “Augmenting Urban Scene Datasets using Diffusion Models,” in *WACV*, 2023.
- [11] Y. Zhang *et al.*, “Large-Scale Pre-training for Image Generation,” in *ICCV*, 2023.
- [12] A. Dosovitskiy *et al.*, “An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale,” in *ICLR*, 2021.
- [13] P. Esser *et al.*, “Taming Transformers for High-Resolution Image Synthesis,” in *CVPR*, 2021.

APPENDIX

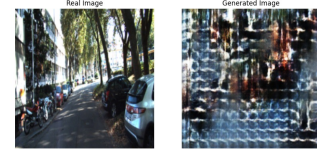
In this appendix, we present additional model output images demonstrating the qualitative results of each generative approach. These images are resized for compactness.



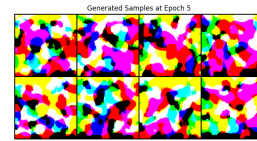
(a) (a) Conditional GAN



(b) (b) U-Net Diffusion



(c) (c) U-Net ViT



(d) (d) ResNet based U-Net



(e) (e) Pre-trained Diffusion

Fig. 4. Additional samples of outputs from each model, resized for clarity: (a) Conditional GAN, (b) U-Net Diffusion, (c) U-Net ViT, (d) ResNet based U-Net, and (e) Pre-trained Diffusion.