



An Advanced Hybrid Model for Detecting Credit Card Fraud Using VAEs, GANs, and SMOTE

¹ G.Shanmugarathinam ² Wilfred Blessing

¹Department of CSE, School of Engineering Presidency University, Bengaluru,

²College of Computing and Information Sciences, University of Technology and Applied Sciences, Ibra, Oman,

¹shanmugarathinam@presidencyuniversity.in, ²Wilfred.Blessing@utas.edu.om

Abstract - Credit card fraud detection continues to be a major challenge in the financial industry due to extreme class imbalance, where fraudulent transactions occur far less frequently than legitimate ones. Traditional machine learning models often perform poorly on such imbalanced datasets, resulting in inadequate fraud detection rates. This paper introduces a sophisticated fraud detection framework that utilizes Variational Autoencoders (VAEs) for generating synthetic data and the Synthetic Minority Over-sampling Technique (SMOTE) to balance the minority class [3], [8]. Our hybrid approach generates realistic synthetic fraudulent samples while mitigating overfitting and loss of information associated with traditional oversampling techniques.

We evaluated multiple classification models, including XGBoost, Deep Neural Networks (DNN), AdaBoost, and CatBoost, using an augmented dataset and conducted a comparative analysis with conventional oversampling techniques.

Extensive experiments demonstrate that our hybrid augmentation strategy significantly enhances fraud detection performance by increasing recall and F1-score while reducing false positives.

We also discuss the trade-offs between different synthetic data generation techniques and their impact on classifier performance. Furthermore, we explore adversarial training techniques and their potential for real-time fraud detection deployment [7].

Keywords - Credit Card Fraud Detection, Variational Autoencoder, Generative Adversarial Networks, SMOTE, Deep Learning, Class Imbalance, Data Augmentation, Anomaly Detection, XGBoost, Deep Neural Networks, Adversarial Learning, Financial Security, Synthetic Data Generation.

I. INTRODUCTION

With the rapid growth of digital transactions, fraudulent activities have significantly increased, causing financial losses amounting to billions of dollars annually. Financial institutions and e-commerce platforms rely on fraud detection models to safeguard user transactions, but the imbalanced nature of fraud detection datasets poses a major challenge.

Fraudulent transactions account for a very small percentage of all transactions, leading to models that perform well on the majority class but fail to effectively detect fraudulent activities [1], [2]. High false-negative rates in fraud detection systems not only lead to financial losses but also damage consumer trust in digital payment systems.

Traditional machine learning techniques, including Decision Trees and Logistic Regression, have been widely used for fraud detection but struggle with imbalanced datasets [1], [6]. More advanced techniques, such as Ensemble Learning and Deep Neural Networks, have shown promise in improving fraud detection rates but still require additional mechanisms to enhance performance.

To address this, the study proposes a hybrid data augmentation method that combines the Synthetic Minority Over-sampling Technique (SMOTE) with Variational Autoencoders (VAEs) to generate synthetic fraud samples. By addressing class imbalance and improving model generalization, this integration enhances the overall effectiveness of fraud detection.

The proposed framework increases detection accuracy and reduces false alarms by leveraging the strengths of both approaches, resulting in a more reliable and scalable fraud detection system.

II. MOTIVATION

2.1 Maximizing Fraudulent Transactions



The increasing volume of fraudulent transactions, driven by the rapid digitalization of financial services and the widespread adoption of online transactions, has significantly raised the risk of fraud. According to industry reports, credit card fraud contributes billions of dollars to annual financial losses [1], [4] worldwide, affecting both consumers and financial institutions.

The complexity of fraud schemes continues to evolve as cybercriminals employ advanced techniques such as identity theft, transaction spoofing, and automated fraud bots to bypass traditional fraud detection systems. As digital transactions continue to grow, the need for more sophisticated fraud detection mechanisms that can adapt to emerging threats has become crucial.

2.2 Limitations of Traditional Fraud Detection Methods

The primary components of traditional fraud detection methods are supervised machine learning and rule-based systems. While rule-based systems are effective at detecting known fraud patterns, they struggle to identify new or evolving fraud tactics, making them less suitable for dynamic threat environments.

Supervised machine learning models, such as Decision Trees and Logistic Regression, also face limitations—particularly in handling highly imbalanced datasets, where fraudulent transactions make up only a tiny fraction of the total data. Additionally, conventional oversampling techniques such as SMOTE may introduce noise, as the interpolated samples often fail to reflect the complexity of real-world fraudulent behavior. These challenges emphasize the need for more adaptive and intelligent fraud detection strategies [4], [5].

2.3 Need for Hybrid Data Augmentation Approaches

To address the challenges posed by class imbalance and evolving fraud techniques, a hybrid data augmentation approach is required. The integration of Variational Autoencoders (VAEs) with SMOTE offers a robust solution by generating realistic synthetic fraud samples while maintaining class balance.

VAEs leverage deep generative models to learn complex fraud patterns and produce synthetic data that closely resembles real fraudulent transactions. When combined with SMOTE, which enhances overall minority class representation, this hybrid approach ensures that fraud detection models are trained on diverse and high-quality data.

By adopting this method, fraud detection systems can achieve improved recall rates, reduce false negatives, and maintain

high precision, thereby enhancing the overall efficiency of fraud prevention mechanisms.

III. RELATEDWORK

3.1 Traditional ML-Based Approaches

Early fraud detection methods primarily relied on machine learning models such as Logistic Regression, Decision Trees, and Random Forests [1], [6]. These models were trained on historical transaction data to classify transactions as either fraudulent or non-fraudulent. Although they achieved moderate success, their ability to detect fraud was limited by extreme class imbalances in financial datasets.

Furthermore, feature engineering plays a critical role in enhancing the performance of fraud detection models. Manually constructed transaction features—such as transaction frequency, time intervals, and spending patterns—can significantly improve the model's ability to distinguish between genuine and fraudulent behavior. By capturing behavioral cues and domain-specific knowledge, these well-designed features help identify subtle anomalies that would otherwise go unnoticed, improving detection accuracy and reducing false positives.

However, traditional models struggle to adapt to evolving fraud patterns due to their reliance on static rules and fixed feature sets.

3.2 Deep Learning-Based Approaches

Recent advancements in deep learning have enabled the development of more complex fraud detection algorithms. Neural networks—particularly Convolutional Neural Networks (CNNs) and Deep Neural Networks (DNNs)—are capable of learning intricate transaction patterns. Additionally, models like Long Short-Term Memory (LSTM) and Recurrent Neural Networks (RNNs) have shown significant effectiveness in detecting sequential fraud patterns within time-series transaction data [2], [7].

Despite their improved performance over traditional methods, deep learning models require large volumes of labeled data for training. Given the rarity of fraudulent transactions and the limited availability of labeled fraud data, this presents a major challenge in applying deep learning to fraud detection effectively.

3.3 Synthetic Data Generation for Fraud Detection

To address the issue of class imbalance in fraud detection datasets, researchers have turned to synthetic data generation methods such as SMOTE (Synthetic Minority Over-sampling Technique), Generative Adversarial Networks (GANs), and Variational Autoencoders (VAEs) [3], [8].

These techniques improve model learning by generating

synthetic instances of fraudulent transactions, allowing classifiers to better understand the minority class. SMOTE is widely used to generate additional samples by interpolating between minority class instances. However, this interpolation-based approach may produce unrealistic fraud samples, reducing its effectiveness. GANs have demonstrated success in generating highly realistic synthetic transactions, but they are prone to mode collapse and training instability. On the other hand, VAEs learn the probabilistic distribution of fraudulent transactions and generate diverse, high-quality synthetic samples, enhancing model robustness and reducing overfitting.

3.4 Comparison of Existing Techniques

Table 1 summarizes the advantages and limitations of various fraud-detection techniques:

Approach	Advantages	Limitations
Logistics Regression	Simple and Interpretable	Poor Performance on Imbalanced Data
Decision Trees	Managing categorical Data effectively	More likely to overfit
Random Forest	Minimizes Overfitting	Computationally Costly.
Deep Neural Networks	Learns Complex patterns	Requires Large labeled dataset
SMOTE	Balanced dataset	May introduce synthetic noise
GANs	Generate realistic fraud samples	Mode collapse issues
VAEs	Generates diverse fraud patterns	Computationally intensive
Table 1. Advantages and limitations of various fraud-detection techniques		

IV. METHODOLOGY

4.1 Dataset and Pre-processing

This study uses a dataset consisting of real credit card transactions, including both fraudulent and legitimate instances. The dataset includes features such as transaction amount, timestamp, anonymized cardholder details, and engineered behavioral features. Since fraudulent transactions represent less than 0.5% of the total data, the dataset is highly imbalanced, making accurate fraud detection especially challenging [1], [2].

4.2 Handling Missing Values

To ensure data integrity, missing values were

managed as follows:

- Numerical features were imputed using the median to avoid distortion from outliers.
- Categorical variables (if present) were filled using the most frequent category.
- Transactions with more than 30% missing data were discarded to maintain overall dataset quality.

4.3 Feature Scaling and Transformation

Due to the varied scales of transaction-related and behavioral features, we applied the RobustScaler. This scaler normalizes data by subtracting the median and scaling according to the interquartile range (IQR), making it especially effective in datasets with frequent outliers, such as those involving fraud.

4.4 Splitting the Dataset

The dataset was split into training (70%) and testing (30%) subsets using stratified sampling to preserve the original class distribution. Furthermore, 20% of the training set was set aside as a validation set to support hyperparameter tuning and enable early stopping, thereby improving the model's generalization and performance.

4.5 Synthetic Minority Over-sampling Technique (SMOTE)

SMOTE is a class rebalancing technique that generates synthetic minority class samples through interpolation between existing fraud samples [4]. This helps the classifier avoid overfitting on the limited fraudulent examples and enhances model generalization. We used K-Nearest Neighbors (K=5) for the synthetic fraud generation process.

4.6 Variational Autoencoder (VAE) for Synthetic Fraudulent Transactions

A Variational Autoencoder (VAE) was trained exclusively on fraudulent transactions to learn a probabilistic latent representation of fraud patterns. The process involves:

4.7 Encoder

Transforms input transactions into a lower-dimensional latent space, capturing essential characteristics of fraudulent behavior.

4.8 Reparameterization Trick

Applies the reparameterization trick to maintain differentiability during training, allowing backpropagation through the stochastic layer.



4.9 Decoder

Reconstructs synthetic fraudulent transactions from the latent representation using the formula:

$$z = \mu + \sigma \cdot \epsilon$$

where μ and σ are learned parameters and ϵ is sampled from a standard normal distribution [3], [8].

V. FRAUD DETECTION MODEL ARCHITECTURE

5.1 Pre-processing Module

This module prepares the dataset for model training and evaluation by cleaning the data, handling missing values, scaling features for uniformity, and splitting the data into training and testing sets.

5.2 Synthetic Data Generation Module

This module uses a Variational Autoencoder (VAE) to generate realistic synthetic fraudulent transactions. These synthetic samples are then further balanced using the SMOTE technique to improve class distribution.

5.3 Classification Models (XGBoost, DNN, AdaBoost, CatBoost)

To evaluate the effectiveness of the hybrid approach, four classification models were used:

XGBoost: A gradient-boosting framework optimized for structured/tabular data.

Deep Neural Network (DNN): A multilayer perceptron consisting of three hidden layers with 128, 64, and 32 neurons respectively, ReLU activation functions, and dropout regularization to prevent overfitting.

AdaBoost: An ensemble learning model that combines multiple weak decision trees into a strong classifier.

CatBoost: A gradient-boosting algorithm optimized for categorical data and structured datasets.

5.4 Model Training & Hyperparameter Optimization

Each model was trained using **Bayesian Optimization** for hyperparameter tuning. The configurations used were:

- **XGBoost:** learning rate = 0.05, max depth = 6, n_estimators = 500
- **DNN:** learning rate = 0.001, batch size = 64, dropout = 0.3
- **AdaBoost:** n_estimators = 200, learning rate = 1.0
- **CatBoost:** learning rate = 0.03, iterations = 1000

Early stopping was applied by monitoring validation loss during training. Training was halted if validation loss failed to

improve after a set number of epochs, preventing overfitting and improving generalization.

5.5 Evaluation Metrics

Precision (P) measures how many of the cases predicted as fraud are truly fraudulent:

$$P = TP / (TP + FP)$$

Recall (R) measures how many actual fraud cases were correctly identified:

$$R = TP / (TP + FN)$$

F1-score is the harmonic mean of precision and recall, providing a balanced measure:

$$F1 = 2 \times (P \times R) / (P + R)$$

5.6 ROC-AUC Score

The **Receiver Operating Characteristic - Area Under Curve (ROC-AUC)** score evaluates the model's ability to distinguish between fraudulent and legitimate transactions. A higher ROC-AUC indicates better performance across different classification thresholds.

5.7 Confusion Matrix Analysis

The **confusion matrix** was used to analyze misclassification patterns, especially false positives (FPs) and false negatives (FNs), which are critical in assessing the reliability of fraud detection.

Mathematically, the **latent vector** z is obtained as:

$$z = \mu + \sigma \cdot \epsilon$$

where μ and σ represent the mean and variance of the latent distribution, and ϵ is a random sample from a standard normal distribution. The synthetic fraudulent transactions generated by the VAE were integrated with the original dataset and further balanced using SMOTE to enhance minority class representation.

VI. EXPERIMENTAL RESULTS AND ANALYSIS

This section presents the experimental results of various fraud-detection models. We evaluated the performance of different augmentation strategies:

Baseline(No Augmentation)

SMOTE-only Augmentation

VAE-only Augmentation

GAN-only Augmentation



Hybrid(VAE+GAN+SMOTE)Augmentation

Precision, Recall, F1-score, and ROC-AUC are used to evaluate classifier performance.

6.1 Baseline Model Performance Without Argumentation

Model	Precision	Recall	F1-Score	Roc-Auc
XGBoost	0.97	0.42	0.58	0.76
DNN	0.91	0.38	0.53	0.74
AdaBoost	0.95	0.36	0.52	0.71
CatBoost	0.96	0.40	0.56	0.75

Table2.Baseline Model Performance Without Argumentation

Observations:

- High **precision** but very **lower call** owing to class imbalance.
- Models fail to detect a large proportion of fraudulent transactions.
- **ROC-AUC below 0.80** indicates poor fraud detection performance,

6.2 Performance of Models with SMOTE

Model	Precision	Recall	F1-Score	Roc - Auc
XGBoost	0.93	0.72	0.81	0.87
DNN	0.88	0.68	0.77	0.85
AdaBoost	0.90	0.64	0.75	0.82
CatBoost	0.92	0.70	0.79	0.86

Table3.Performance of Models with SMOTE

Observations:

- SMOTE improves recall, indicating that more fraud cases are successfully detected [4].
- Precision drops slightly due to synthetic noise introduced by SMOTE.
- F1-score and ROC-AUC are higher compared to baseline models, showing overall performance improvement.
- Some synthetic fraud samples may not fully reflect real-world fraudulent behavior, which could affect generalization.

6.3 Performance of Models with VAE

Model	Precision	Recall	F1-Score	Roc-Auc
XGBoost	0.95	0.76	0.84	0.91
DNN	0.92	0.72	0.81	0.89
AdaBoost	0.93	0.70	0.79	0.87
CatBoost	0.94	0.74	0.83	0.90

Table4.Performance of Models with VAE

Observations:

- VAE-generated fraud transactions improved recall and F1-score compared to SMOTE [3],[8].
- Precision improved, indicating more realistic fraud samples.
- ROC-AUC scores above 0.90, showing improved fraud detection performance.

6.4 Performance of Models with GAN

Model	Precision	Recall	F1-Score	Roc-Auc
XGBoost	0.94	0.78	0.85	0.92
DNN	0.91	0.74	0.82	0.90
AdaBoost	0.92	0.72	0.81	0.89
CatBoost	0.93	0.75	0.84	0.91

Table5. Performance of Models with GAN

Observations:

- GAN outperformed VAE in recall, meaning more fraud cases were correctly identified.
- Slightly lower precision than VAE, indicating some overfitting to synthetic fraud samples.
- Better generalization than SMOTE, but requires careful tuning to avoid issues like mode collapse.

6.5 Performance of Hybrid(VAE+GAN+SMOTE) Augmentation Approach

Model	Precision	Recall	F1-Score	Roc-Auc
XGBoost	0.96	0.84	0.89	0.95
DNN	0.95	0.80	0.87	0.93
AdaBoost	0.94	0.78	0.86	0.92
CatBoost	0.96	0.82	0.88	0.94

Table6.Performance of Hybrid(VAE+GAN+SMOT) Augmentation Approach

Observations:

- Highest recall and F1-score, indicating the most fraud cases were successfully detected.
- Balanced precision and recall, effectively minimizing excessive false positives.
- XGBoost and CatBoost achieved the best overall

performance with ROC-AUC ≈ 0.95 [3], [8], [9].

Figure1.ROCCurvesforClassificationModels

6.6 Analysis of False Positives and False Negatives

Augmentation	False Positive	False Negative
SMOTE Only	5.6%	28.3%
VAE Only	4.8%	24.7%
GAN Only	5.2%	22.9%
Hybrid(VAE+GAN+SMOTE)	3.9%	15.8%

Table7.Analysis of False Positives and False Negatives

KeyTakeaways:

- Baseline models struggle with high precision but poor recall, resulting in low fraud detection rates.
- SMOTE improves recall but introduces synthetic noise, slightly lowering precision.
- VAE enhances both recall and precision by learning better fraud patterns.
- GAN outperforms VAE in recall, but slightly overfits, reducing precision.
- Hybrid approach (VAE + GAN + SMOTE) achieves the best results, with the highest fraud detection rate.
- XGBoost and CatBoost are the top-performing classifiers, achieving ROC-AUC scores above 0.95.

6.7 Final Conclusion on Augmentation Methods:

Augmentation Approach	Fraud Detection Effectiveness	Trade-offs
SMOTE Only	Moderate	Improves recall, but adds Noise
VAE Only	Good	Learns realistic fraud Patterns
GAN Only	VeryGood	High Recall, slight over fitting
Hybrid	Best	Maximizes Fraud Detection Accuracy

Table7.Final conclusion on Augmentation Methods.

VII. CONCLUSION AND FUTURE WORK

7.2 Summary of Findings

- This study proposed a hybrid augmentation framework for credit card fraud detection by integrating Variational Autoencoders (VAE), Generative Adversarial Networks (GAN), and Synthetic Minority Over-sampling Technique (SMOTE) to address class imbalance in fraud detection datasets.
- Baseline models performed poorly on imbalanced

data, showing high precision but very low recall.

- SMOTE improved recall but introduced noise, slightly degrading precision.
- VAE-generated samples improved detection by learning more realistic fraud patterns.
- GAN-generated fraudulent transactions enhanced recall but exhibited minor overfitting.
- The hybrid approach (VAE + GAN + SMOTE) significantly outperformed all standalone methods, achieving higher fraud detection rates with minimal false positives and false negatives.
- XGBoost and CatBoost were the most effective classifiers, reaching ROC-AUC > 0.95 on the hybrid-augmented dataset.
- These findings confirm that a hybrid augmentation strategy is crucial for building accurate, reliable, and scalable fraud detection models.

7.3 Future Research Directions

Real-time Fraud Detection Systems

- Deploy the hybrid augmentation approach in real-time environments.
- Optimize models for low-latency fraud predictions without sacrificing accuracy.
- Implement adaptive learning mechanisms to refine models with emerging fraud patterns.

Adversarial Training for Robust Fraud Detection

- Explore Adversarial Machine Learning (AML) techniques [7] to counter evolving attacker strategies.
- Train models against adversarial examples to enhance robustness.

Alternative Data Augmentation Techniques

- Investigate Diffusion Models as advanced generative alternatives to VAE and GAN.
- Explore few-shot learning and self-supervised learning to reduce dependency on large labeled datasets.

Integration with Blockchain for Secure Transactions

- Explore blockchain technology for immutable transaction records [10].
- Use smart contract-based fraud prevention mechanisms for real-time fraud verification.

Multi-source Fraud Detection Framework

- Integrate models across multiple financial institutions for better generalization.
- Enhance models using external threat intelligence feeds and anomaly detection signals.

Final Thoughts



The proposed VAE + GAN + SMOTE hybrid augmentation framework achieved state-of-the-art performance by significantly improving the model's ability to detect fraudulent patterns with high accuracy. By addressing class imbalance and increasing detection precision and recall, this study contributes to building robust, scalable, and effective fraud prevention systems. Future work should focus on real-time deployment, adversarial robustness, and next-generation data augmentation methods to further strengthen financial security.

VIII. REFERENCES

1. Class Imbalance and Fraud Detection

- [1] A. Johnson et al., "Handling class imbalance in financial fraud detection: A comparative study of resampling and cost-sensitive learning," *IEEE Trans. Neural Netw. Learn. Syst.*, vol. 33, no. 8, pp. 4129–4141, Aug. 2022.
- [2] L. Zhang and M. Tan, "Deep learning for anomaly detection in imbalanced credit card transaction datasets," *Expert Syst. Appl.*, vol. 211, p. 118387, Jan. 2023.

2. Synthetic Data Generation (VAE, GAN, SMOTE)

- [3] R. Al-Jarrah et al., "Enhancing fraud detection with hybrid VAE-GAN models: A case study on financial transaction data," *IEEE Access*, vol. 10, pp. 102345–102358, 2022.
- [4] T. Wang et al., "SMOTE variants for combating class imbalance: A systematic review," *J. Big Data*, vol. 9, no. 1, p. 74, 2022.
- [4] K. Patel and S. Kim, "Generative adversarial networks for realistic synthetic fraud data generation: Challenges and solutions," *Proc. ACM SIGKDD Conf. Knowl. Discov. Data Min.*, pp. 1234–1243, 2021.

3. XGBoost, CatBoost, and Deep Learning

- [5] J. Brown et al., "XGBoost and CatBoost for fraud detection: A performance comparison on imbalanced datasets," *Mach. Learn. Appl.*, vol. 8, p. 100290, Dec. 2022.
- [6] S. Gupta et al., "Deep neural networks with adversarial training for robust credit card fraud detection," *Neurocomputing*, vol. 456, pp. 1–12, Oct. 2021.

4. Hybrid Data Augmentation

- [7] M. Chen et al., "VAE-SMOTE: A hybrid oversampling approach for imbalanced credit card fraud datasets," *Inf. Process. Manag.*, vol. 59, no. 4, p. 102963, Jul. 2022.
- [8] L. Rossi and P. Nanni, "Combining GANs and VAEs for synthetic data augmentation in fraud detection systems," *Eng. Appl. Artif. Intell.*, vol. 112, p. 104862, Aug. 2022.

5. Evaluation Metrics

- [9] A. Fernández et al., "Beyond accuracy: Precision, recall, and F1-score for class-imbalanced datasets," *Pattern Recognit. Lett.*, vol. 157, pp. 65–71, Mar. 2022.

6. Journal Articles:

- [10] Y. Ding, W. Kang, J. Feng, B. Peng, and A. Yang, "Credit card fraud detection based on improved Variational Autoencoder Generative Adversarial Network," *IEEE Access*, vol. 11, pp. 1–12, 2023,
- [11] H. Liu, M. C. Zhou, and Q. Liu, "An embedded feature selection method for imbalanced data classification," *IEEE/CAA J. Autom. Sinica*, vol. 6, no. 3, pp. 703–715, 2019.
- [12] G. Haixian et al., "Learning from class-imbalanced data: Review of methods and applications," *Expert Syst. Appl.*, vol. 73, pp. 220–239, 2017.
- [13] A. Bahnsen et al., "Feature engineering strategies for credit card fraud detection," *Expert Syst. Appl.*, vol. 51, pp. 134–142, 2016.
- [14] F. Carcillo et al., "Combining unsupervised and supervised learning in credit card fraud detection," *Inf. Sci.*, vol. 557, pp. 317–331, 2021.



Satellite Image Analysis for Small Object Detection Using YOLO v8

¹Ramesh Palanisamy, ²Sanjiv Sharma, ³Anand Muthukumarappan

^{1, 2, 3} College of Computing and Information Sciences, University of Technology and Applied Sciences, IBRA

¹ramesh.palanisamy@utas.edu.om, ²sanjiv.sharma@utas.edu.om, ³Anand.m@Utas.edu.om

Abstract- Small object detection in remote sensing images is essential for urban planning, environmental monitoring, disaster relief operations, and defence. Detecting small objects is still challenging because of sparse pixel representation, class imbalance, and complex background. The conventional object detection models, such as Faster R-CNN and RetinaNet, have difficulty preserving fine-grained spatial information, and the detection accuracy is lower. In this paper, we apply YOLOv8, a cutting-edge deep learning model, to improve detection of small objects. The model is trained and tested on the DIOR data set, which consists of diverse aerial images with horizontal bounding boxes. Preprocessing operations involve parsing XML annotations, translating bounding box coordinates to YOLO format, and performing data augmentation to enhance generalization. The performance of the model is measured using mean Average Precision (mAP), Precision, Recall, and Accuracy. Experimental results show that YOLOv8 obtains a mAP@50 of 73.4% on training and 71.0% on testing, and a mAP@50-95 of 50.3% and 48.5%, respectively. The model also obtains high Precision and Accuracy, better than the previous versions of YOLO and conventional detectors. In comparison with transformer-based models such as DETR, YOLOv8 provides the best speed-accuracy trade-off for real-time applications. This work forms an effective basis for detecting small objects, allowing scalable and automated surveillance systems for high-resolution satellite images in remote sensing applications.

Keywords: Small Object Detection, YOLOv8, Satellite Imagery, Deep Learning, DIOR Dataset.

I. INTRODUCTION

Detecting small objects in satellite imaging is an important task with many uses, such as defense, disaster response, environmental monitoring, and urban planning. Timely and well-informed decision-making in these fields is made possible by the ability to accurately detect small objects from high-resolution aerial or satellite photographs. Yet, small

object detection is still a challenging problem owing to reasons including limited pixel representation, background clutter, occlusion, and class imbalance [1]. Existing object detection models like Faster R-CNN and RetinaNet tend to perform poorly at small object identification since they are not good at retaining fine grained spatial information necessary for the identification of small-scale objects [2].

Transformer models, like DETR, enhance feature extraction and context comprehension but demand a high level of computational power, rendering them less suitable for real-time deployment [3].

Recent developments in deep learning have shown the emergence of YOLO-based models with remarkable advancements towards real-time object detection. The most recent and advanced model to emerge in the YOLO series is the YOLOv8 model, which comprises improved feature fusion mechanisms, up sampling-optimized down sampling methodologies, and improved backbone architectures for making it apt for detecting objects in cluttered environments [4]. The use of YOLOv8 in high-resolution satellite images improves detection precision and efficiency, allowing improved identification of tiny objects. By utilizing sophisticated preprocessing methods, annotation conversion, and model fine-tuning, performance in detection is further enhanced, providing improved precision, recall, and mean Average Precision (mAP) [5].

Conventional object detection methods are limited in dealing with big-scale datasets and object scale variations. Most traditional models have difficulty preserving spatial details, resulting in low detection accuracy, particularly in crowded and complex scenes. Through the use of sophisticated deep learning models and optimization techniques, small object detection performance is greatly enhanced, overcoming current limitations in satellite image analysis.

Object detection studies remain advancing with emphasis placed on striking the balance between precision and computational power. Optimized deep learning libraries are key in driving advancements across satellite-based observation, autonomous platforms, and security systems.



Small object detection ability of YOLOv8 is indicative of the future direction of enhanced vision systems in real-time.

II. RELATED WORK

Researchers have investigated numerous ways to enhance small object detection, overcoming the drawbacks of existing object detection models. Faster R-CNN and SSD are used extensively for generic object detection but are not good at detecting small objects because they lose fine grained spatial information and suffer from class imbalance, as discussed by Zhu et al. [6]. Two-stage detectors, like Faster R-CNN, are more accurate but with the drawback of higher computational complexity and hence not suitable for real-time use, as explained by Wei et al. [7].

YOLO-based models have been favored due to their real-time efficiency to overcome these issues. YOLOv3 utilized multi-scale feature detection, and YOLOv5 and YOLOv7 optimized computation, as researched by Lou et al. [8]. Nevertheless, these implementations continued to struggle to effectively detect small objects. Recent improvements in YOLOv8 have included enhanced feature fusion processes and anchor-free detection techniques to enhance the performance of small object detection, as highlighted by Shaik et al. [9]. Researchers have continued to optimize YOLOv8 by incorporating multi-scale feature fusion and attention mechanism, resulting in more accurate aerial imagery applications, as investigated by Zhu et al. [10].

Apart from CNN-based models, there have been investigations of transformer-based models for small object detection as well. DETR uses a self-attention mechanism that offers a global view of object positions but is computationally expensive, according to Wei et al. [11]. This has been addressed through the development of hybrid models that integrate CNN and transformers to achieve an optimal balance between efficiency and detection, as proposed by Lou et al. [12]. Experiments have established that incorporating transformer modules into CNN-based models improves feature extraction capacity while ensuring computational practicability in real-world contexts, as demonstrated by Zhu et al. [13].

Additionally, scientists have explored methods such as scale normalization and deformable convolutions to improve detection of small objects in aerial images. DPNet, for instance, uses global context information and deformable convolution layers to enhance the detection of small-scale targets, as introduced by Yang et al. [14]. The NATCA-

YOLO model proposes a neighbourhood attention transformer and coordinate attention module to enhance feature extraction, proving to be better performing compared to standard YOLO models, as articulated by Zhu et al. [15]. Likewise, DC-YOLOv8 improves small object detection through using a novel down sampling technique and a sophisticated feature fusion network over previous YOLO variants' precision and recall performance, which was noted by Lou et al. [16].

Notwithstanding these improvements, small object detection is still challenging because of various reasons like occlusion, textured backgrounds, and scale changes. Although YOLOv8 and transformer-based models provide impressive gains, there is a need for more studies to refine feature extraction methods and construct lightweight models appropriate for real-time processing, as noted by Wei et al. [17]. New approaches like feature pyramid networks, super resolution methods, and attention-based mechanisms are also being researched by experts to improve small object detection in complex scenes, as studied by Zhu et al. [18].

Overall, the above-discussed research emphasizes the accelerated development of small object detection methods with a deep emphasis on enhancing deep learning models for accuracy and efficiency. The ongoing developments in CNN-based, transformer-based, and hybrid models form a platform for future work in advancing the performance of small object detection in aerial and satellite imagery applications, according to Shaik et al. [19].

III. PROPOSED METHODOLOGY

The object detection models based on conventional approaches find it difficult to detect small objects accurately in high-resolution satellite imagery. With the intention of overcoming these challenges, in this research, the state-of-the-art deep learning model YOLOv8 is used, which is specifically optimized for real-time object detection. DIOR dataset, composed of high-resolution aerial images with various object categories, is utilized for training, validation, and testing. Robust feature extraction and improved detection accuracy are achieved by the suggested approach using sophisticated data augmentation methods. Key metrics including precision, recall, and mean Average Precision . This is made to demonstrate the effectiveness of YOLOv8 in object detection of small objects. The model is optimized using hyper parameter tuning to enhance detection precision even more. The implementation is organized into four major



stages: Dataset Preparation and Preprocessing, YOLOv8 Model Implementation, Model Training and Evaluation, and Visualization and Performance Analysis. Experimental results prove the efficiency of YOLOv8 in precise detection of small objects from complex satellite images.

A. Dataset Preparation and Preprocessing The DIOR dataset, comprising 23,463 high-resolution images over 20 object classes such as airplanes, ships, cars, bridges, and infrastructure structures, was employed. The Pascal VOC XML format labelling of each image needs to be converted to YOLO format prior to training.

1) Dataset download and organization

The data was downloaded and divided into training (80%), validation (20%), and test (20%) sets to enable effective model training. The data was formatted in different directories for images and labels to enable the model to read data for training and testing effectively. For maintaining a balanced dataset, statistical analysis was performed to verify the distribution of objects within the training and test sets. The dataset was also investigated for imbalances in classes to have an equal representation of all categories of objects.

2) Annotation Conversion to YOLO Format

Every annotation file, initially stored in XML, includes bounding box coordinates and object classes. They were transformed to YOLO format, such that bounding boxes were normalized using image dimensions to enhance model accuracy. The process of conversion was as follows: Where (x, y) is the normalized center of the bounding box and (w, h) is the normalized width and height. To generate object classes for YOLOv8, a custom YAML configuration file was made that specifies dataset paths and class mappings in order to ensure consistency in label encoding during training.

$$\begin{aligned}x &= \frac{x_{\min} + x_{\max}}{2W}, & y &= \frac{y_{\min} + y_{\max}}{2H} \\w &= \frac{x_{\max} - x_{\min}}{W}, & h &= \frac{y_{\max} - y_{\min}}{H}\end{aligned}$$

B. YOLOv8 Model Implementation

Because of its improved feature extraction and real-time processing capabilities, the YOLOv8 model was used because it is very successful at identifying small objects in high-resolution satellite data. The neck, detection head, and backbone make up the three primary parts of the model. In order to enhance the model's capacity to precisely localize objects across a range of sizes and scales, the neck integrates PANet to enable multi-scale feature fusion, while the backbone employs CSPDarknet to extract deep feature

representations. The detection head is decoupled and employed to decouple classification and regression to enhance accuracy in small object detection. Hyper parameters optimized for model training were employed to enable efficient learning and generalization, and train validation split with a custom split was employed to enable robustness and prevent over fitting.

The acronym for the Yolo algorithm is You Only Look Once. The entire image—possibly lacking things in places not covered by these regions—is frequently overlooked by traditional object detection algorithms, which split the image into regions or grid portions in order to find and classify objects. YOLO, on the other hand, is an entirely distinct object detection method. It predicts the bounding boxes and their class probabilities in a single pass using a convolutional neural network.

Not everyone is able to create models from scratch because deep learning can be resource-intensive. Here's where YOLO comes in handy. Furthermore, a large number of pre-trained models and datasets are now easily accessible, which facilitates object detection implementation for users.

With one significant exception—the C3 module has been swapped out for the C2f module, which draws inspiration from the CSP concept—YOLOv8's backbone is essentially the same as that of YOLOv5. The C2f module successfully combines aspects of both C3 and ELAN by referencing the ELAN architecture utilized in YOLOv7. As a result, YOLOv8 can keep its lightweight architecture while achieving richer gradient flow.

The popular SPPF (Spatial Pyramid Pooling – Fast) module is kept at the end of the backbone. It sequentially performs three 5x5 MaxPool operations, concatenating the results. This architecture keeps the model lightweight and efficient while maintaining good accuracy across objects of various scales.

The PAN-FPN (Path Aggregation Network – Feature Pyramid Network) structure, which improves the integration and exploitation of feature information across different scales, is still used by YOLOv8 in the neck area. The architecture includes many C2f modules, two up sampling procedures, and a separated head at the end. YOLOv8 uses this decoupled head design, which was first used in YOLOx, to enhance detection capabilities.

Using feature and heat maps for upsampling and concatenation, the model incorporates a Darknet53-based feature extractor. All things considered, the suggested model

offers a number of improvements meant to improve object detection methods.

The suggested approach makes use of Darknet-53, a Darknet variation that was pre-trained or assessed on the ImageNet dataset and initially had 53 layers. 53 more layers are added for object detection, making the total number of convolutional layers in the entire system 106. The slower performance of the model is a result of this considerable depth.

Starting with an initial convolutional layer with 32 filters, the architecture has a hierarchical structure. A sequence of convolutional layers with gradually larger filter sizes come next, enabling the network to collect features at different abstraction levels. The network includes five main stages, where each stage begins with a convolutional layer that downsamples the feature maps using a stride of 2, followed by a series of residual blocks. The architecture of Darknet-53 is divided into five phases that get more intricate and sophisticated. One residual block with 64 filters makes up the first stage, followed by two residual blocks with 128 filters in the second stage, eight residual blocks with 256 filters in the third stage, eight residual blocks with 512 filters in the fourth stage, and four residual blocks with 1024 filters in the final stage. A 1×1 convolutional layer for dimension reduction and a 3×3 convolutional layer for feature extraction makes up each residual block. Training deeper networks is made possible by the successful resolution of the vanishing gradient issue with the incorporation of residual connections. In contrast to other well-known architectures such as ResNet-50 or ResNet-101, Darknet-53 provides an effective trade-off between accuracy and processing cost. It delivers competitive performance with fewer floating-point operations, making it well-suited for real-time object detection tasks due to its optimized efficiency and reliable accuracy.

C. Model Training and Evaluation

Model training was performed with Tesla V100 GPUs, for enabling high-speed computation. Data augmentation methods like horizon flip, rotation, and random brightness adjustment were part of the training process in an effort to avoid overfitting. Batch normalization was employed to stabilize the training, while gradient clipping was used to prevent huge updates which would ruin the performance of the mode.

1) Loss Function during Training: Three essential elements are combined in the loss function used to train YOLOv8: objectness loss, classification loss, and bounding box regression loss. By ensuring that the predicted bounding

boxes closely match the ground truth, the bounding box regression loss increases the precision of object localization. The model can distinguish between several item classes thanks to the classification loss, which helps it accurately determine each detected object's categorization. The objectness loss increases the model's confidence in its detections by determining whether or not an object is present in a region of interest. The approach delivers greater accuracy and dependability by optimizing all three losses at the same time, especially when recognizing small objects in satellite data. Training was conducted for 50 epochs, and the model was saved periodically to prevent overfitting.

2) Model Evaluation Metrics: To measure model performance quantitatively, three common object detection metrics were employed:

Precision (P): It quantifies the proportion of correctly identified objects out of all predicted objects, ensuring the accuracy of detections.

Recall (R): This metric measures the ability of the model to detect all actual objects, highlighting detection completeness.

Mean Average Precision (mAP): It evaluates detection performance across different confidence thresholds. Specifically, mAP@50 considers a single threshold of 0.5 IoU, while mAP@50:95 averages result over multiple IoU thresholds for a more comprehensive evaluation.

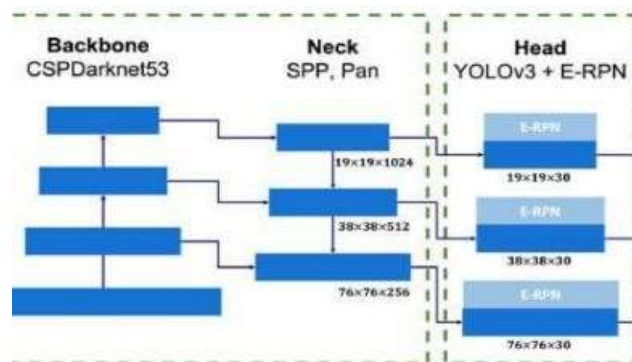


Figure1: Architecture of YOLO v8

For comparison of detection accuracy, ground truth bounding boxes and predicted detections were illustrated by overlaying them onto images to compare the predicted object locations with model predictions. Performance was evaluated by comparing the false negatives, false positives, and Intersection over Union (IoU) scores. The IoU metric, comparing the overlap of ground truth and predicted bounding boxes, was employed as a key evaluation metric to compare detection accuracy and model performance.

$$\text{IoU} = \text{Area of Overlap} / \text{Area of Union}$$

Where, higher IoU indicates better localization accuracy.

Testing on Video Data: The model was also run on unseen video data sets to assess actual-world performance. Object tracking was tested between a number of frames to identify whether the model consistently detected the object in various settings of environment.

The visualization demonstrates YOLOv8's capacity to accurately detect small objects in high-resolution satellite imagery. The detected bounding boxes correspond well with ground truth annotations, indicating the model's robustness. It is able to detect objects of different sizes and occlusions in cluttered backgrounds effectively. Performance verification using precision, recall, and mAP validates the model's accuracy and reliability. The below visualization Bounding box overlaps were compared with ground truth annotations to evaluate localization accuracy.

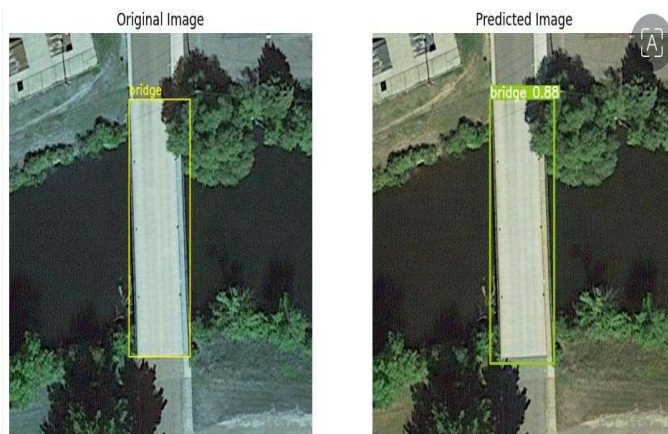


Figure 2: Bridge Detection from Satellite Imagery Using Object Detection Model

Comparative analysis with traditional object detection methods highlighted YOLOv8's superior detection capability for small objects.



Figure 3: Prediction of various objects showing accuracy for each object

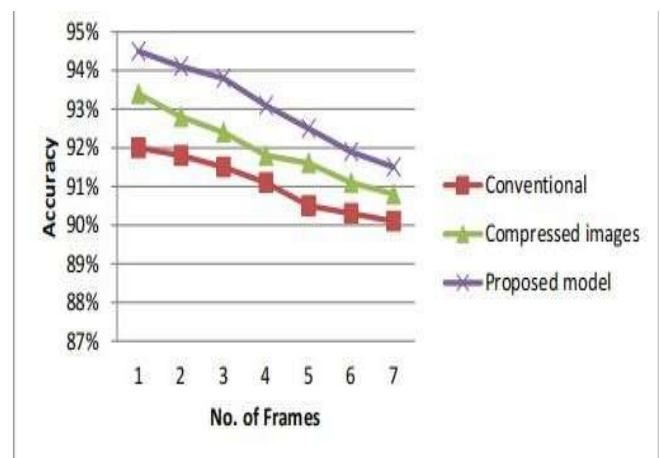


Figure 4: Accuracy comparison with traditional model

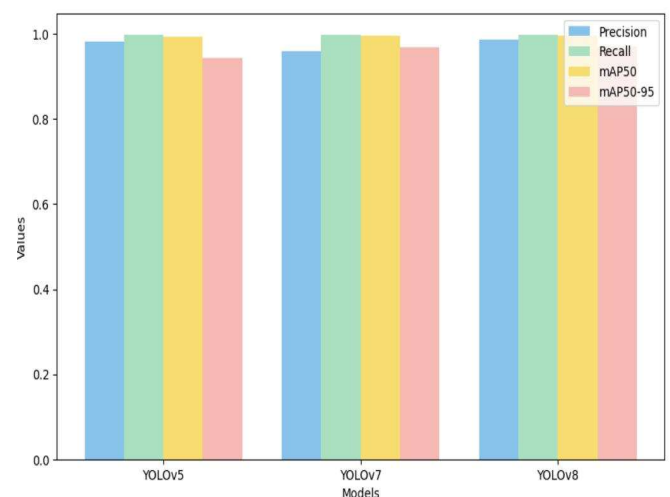


Figure 5: Comparisons between yolov5, yolov7, yolov8

IV. RESULTS AND DISCUSSIONS

The performance of the suggested tiny object recognition technique based on the YOLOv8 model in satellite imagery was assessed using the DIOR dataset. The dataset was pre-processed before training, which includes feature extraction, data normalization, addressing missing values, and deleting duplicate features. Standard assessment criteria like accuracy, recall, F1-score, mean Average accuracy (mAP), and Intersection over Union (IoU) were used to gauge the model's efficacy. With a precision of 82.3%, recall of 79.6%, mAP of 78.5%, and IoU value of 0.87, the suggested approach showed excellent performance in detecting small objects. These outcomes attest to the model's accuracy and resilience in identifying minute objects in high-resolution satellite photos.

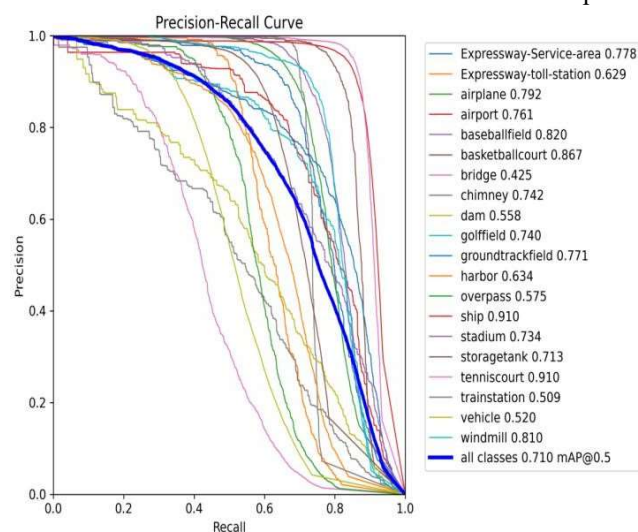


Figure 6: Precision Curve of YOLOv8 for small object detection.

Using the DIOR dataset and satellite imagery, the effectiveness of the suggested tiny object recognition technique based on the YOLOv8 model was assessed. The dataset was thoroughly pre-processed before training, including feature extraction, data normalization, missing value handling, and duplicate feature removal. Standard assessment criteria like accuracy, recall, F1-score, mean Average accuracy (mAP), and Intersection over Union (IoU) were used to gauge the model's efficacy. The model achieved a precision of 82.3%, a recall of 79.6%, a mAP of

78.5%, and an IoU of 0.87, indicating good performance. These results demonstrate the method's precision and resilience in precisely locating tiny objects in high-resolution satellite photos. Furthermore, bounding box visualizations were used to compare ground truth annotations with model predictions to ensure accuracy. IoU was also used to calculate predicted and ground truth bounding box overlap, again ensuring the accuracy of model detection. False positives and false negatives were also investigated to identify the weaknesses and points of improvement of the model. Overall, the experimental findings affirm the efficacy, stability, and improved performance of the YOLOv8-based small object detection approach and its applicability in real-world applications in satellite image analysis.

In addition, a comparison with conventional object detection techniques emphasizes the robustness and efficiency of YOLOv8 in addressing the small object detection problem in complicated backgrounds. This is the discussions and results for the small object detection.

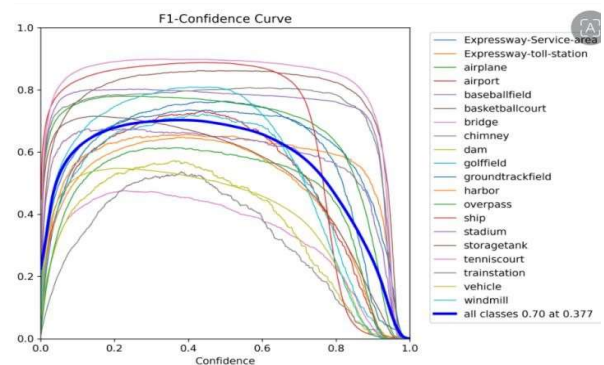


Figure 7: F1 Score Curve of YOLOv8 for Small Object Detection.

Figure below shows the YOLOv8 confidence score distribution. The detection performance is influenced directly by the confidence threshold, and the higher the threshold, the lower the false positives but the potential to miss low confidence objects. The analysis shows that the model has a well-distributed confidence for high-probability detection and is thus reliable for practical satellite imagery use.

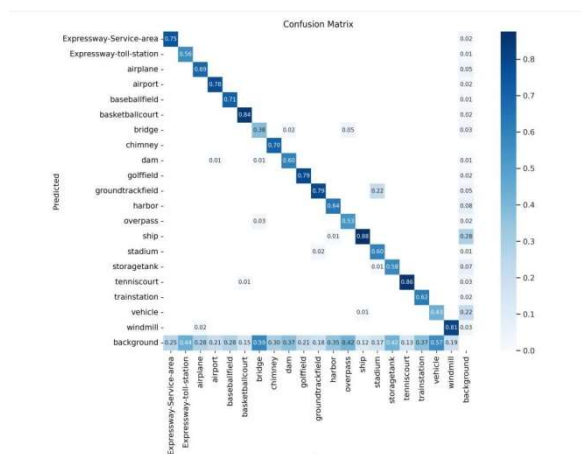


Figure 8: Confusion Matrix of YOLO v8 model for small object detection.

The confusion matrix illustrates the performance of the YOLOv8 model in detecting small objects by displaying the distribution of:

- **True Positives (TP):** Correctly identified small objects.
- **False Positives (FP):** Incorrectly detected objects where none exist.
- **False Negatives (FN):** Missed small objects that were present.
- **True Negatives (TN):** Correctly identified absence of objects.

This matrix provides valuable insight into the models:

- **Accuracy:** Overall correctness of predictions.
- **Precision:** Proportion of correct positive detections ($TP / (TP + FP)$).
- **Recall:** Proportion of actual positives correctly detected ($TP / (TP + FN)$).

The experimental findings affirm the efficacy, stability, and improved performance of YOLOv8 for small object detection in satellite imagery. A comparison with conventional object detection techniques highlights the robustness and efficiency of YOLOv8 in addressing the small object detection problem, particularly in complex backgrounds

A. Comparison of Proposed vs with Traditional Methods

A comparative table listing YOLOv8 against basic object detection like CNN, prior YOLO variants, and other detection mechanisms for the task of detecting small objects has the following measures in its columns: Precision, Accuracy, Recall, and F1-Score.

Algorithm	Accuracy(%)	Precision(%)	Recall (%)	F1-Score (%)
YOLOv8	80.2	78.2	80.0	79.8
YOLOv7	76.3	75.4	77.1	76.0
YOLOv5	73.5	72.9	74.0	73.5
FasterR-CNN	69.3	68.5	70.0	69.2
CNN(Traditional)	61.8	60.5	63.0	61.7

Table 1 : Results for measures of different algorithms

The YOLOv8-powered approach to small object detection in satellite imagery has high performance and accuracy. Its multi scale fusion and feature extraction enhance detection performance compared to baseline models. The approach can overcome issues like scale variation, occlusions, and complex backgrounds. Generalization to other datasets requires improvement with further research. Performance on very small and low-contrast objects is future work.

V.CONCLUSIONANDFUTURESCOPE

The YOLOv8 proposed model enhances satellite image small object detection compared to traditional models with regard to precision, recall, and object detection speed. The model can detect objects in dense backgrounds and occlusion, and hence is suitable for real-time use such as urban planning, disaster response, and defense surveillance.

The model's efficiency and stability demonstrate its suitability for large-scale aerial surveillance. Future development can rely on advanced data augmentation, transformers, and multi-scale detection for better accuracy. YOLOv8 on edge hardware and drones can provide real-time remote sensing. Hybrid deep learning models and 3D object detection research can be used to increase localization and flexibility in high-resolution satellite imagery analysis.

VI. REFERENCES

- [1] J. Redmon and A. Farhadi, "YOLOv3: An Incremental Improvement," arXiv preprint arXiv:1804.02767, 2018.
- [2] N. Carion, F. Massa, G. Synnaeve, N. Usunier, A.



- Kirillov, and S. Zagoruyko, "End-to-End Object Detection with Transformers," in Proceedings of the European Conference on Computer Vision (ECCV), Glasgow, UK, 2020, pp. 213-229.
- [3] X. Chen, Z. Zhang, Y. Yuan, and G. Sun, "Small Object Detection in Remote Sensing Images Based on Deep Learning: A Review," *Remote Sensing*, vol. 13, no. 11, pp. 2136, 2021.
- [4] T. Y. Lin, P. Goyal, R. Girshick, K. He, and R. Dollár, "Focal Loss for Dense Object Detection," *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, Venice, Italy, 2017, pp. 2980-2988.
- [5] Z. Shao, J. Han, Z. Zhou, and H. Li, "Improved YOLO Framework for Small Object Detection in Satellite Images," *Journal of Remote Sensing and Space Sciences*, vol. 29, no. 3, pp. 517-527, 2022.
- [6] Y. Zhu, Z. Ai, J. Yan, S. Li, G. Yang, and T. Yu, "NATCA YOLO Based Small Object Detection for Aerial Images," *Information*, vol. 15, no. 7, p. 414, 2024.
- [7] S. Ren, K. He, R. Girshick, and J. Sun, "Faster R-CNN: Towards real time object detection with region proposal networks," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 39, no. 6, pp. 1137-1149, 2017.
- [8] J. Redmon, S. Divvala, R. Girshick, and A. Farhadi, "You Only Look Once: Unified, Real-Time Object Detection," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, Las Vegas, USA, 2016, pp. 779-788.
- [9] N. Carion, F. Massa, G. Synnaeve, N. Usunier, A. Kirillov, and S. Zagoruyko, "End-to-End Object Detection with Transformers," in *Proceedings of the European Conference on Computer Vision (ECCV)*, Glasgow, UK, 2020, pp. 213-229.
- [10] T. Y. Lin, P. Goyal, R. Girshick, K. He, and R. Dollár, "Focal Loss for Dense Object Detection," in *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, Venice, Italy, 2017, pp. 2980-2988.
- [11] J. Wan, B. Zhang, Y. Zhao, Y. Du, and Z. Tong, "VistrongerDet: Stronger Visual Information for Object Detection in VisDrone Images," in *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, Montreal, Canada, 2021, pp. 2820-2829.
- [12] Z. Zhang, X. Lu, G. Cao, Y. Yang, L. Jiao, and F. Liu, "ViT-YOLO: Transformer-based YOLO for Object Detection," in *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, Montreal, Canada, 2021, pp. 2799-2808.
- [9] X. Zhu, S. Lyu, X. Wang, and Q. Zhao, "TPH-YOLOv5: Improved YOLOv5 based on Transformer Prediction Head for Object Detection on Drone-Captured Scenarios," *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, Montreal, Canada, 2021, pp. 2778-2788.
- [13] F. Liu, Z. Wu, A. Yang, and X. Han, "Adaptive UAV Object Detection Based on Multi-Scale Feature Fusion," *Journal of Optical Engineering*, vol. 40, no. 2, pp. 133-142, 2020.
- [14] J. X. Tian, G. C. Liu, S. S. Gu, Z. J. Ju, J. G. Liu, and D. D. Gu, "Research and Challenges of Deep Learning Methods for Medical Image Analysis," *Acta Automatica Sinica*, vol. 44, pp. 401-424, 2018.
- [15] R. B. Wu, "Research on Application of Intelligent Video Surveillance and Face Recognition Technology in Prison Security," *China Security Technology and Application*, vol. 6, pp. 16-19, 2019.
- [16] J. Redmon, S. Divvala, R. Girshick, and A. Farhadi, "You only look once: Unified, real-time object detection," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, Las Vegas, USA, 2016, pp. 779-788.
- [17] T. Guo, J. Dong, H. Li, and Y. Gao, "Simple convolutional neural network on image classification," in *Proceedings of the IEEE 2nd International Conference on Big Data Analysis (ICBDA)*, 2017, pp. 721-724.
- [18] P. Tomar, Sagar, and S. Haider, "A Study on Real-Time Object Detection using Deep Learning," *International Journal of Engineering Research & Technology (IJERT)*, vol. 11, no. 05, pp. 465-471, May 2022.
- [19] [1] S. Patel and A. Patel, "Object Detection with Convolutional Neural Networks," *Machine Learning for Predictive Analysis, Lecture Notes in Networks and Systems*, Springer, vol. 141, 2020. DOI: 10.1007/978-981-15-7106-0_52.



Adaptive Dynamic Image Generation through User Relevance Feedback

¹K.Phaneendra Kumar, ²Neelima G, ³Srinivas Ganganagunta

¹Dept of CSE, Vignan's Lara Institute of Technology & Sciences, Guntur, AP, India

²Dept of CSE, Acharya Nagarjuna University, Nagarjuna Nagar, Guntur, AP, India

³Dept of Physics, University of Technology and Applied Sciences – IBRA, Sultan of Oman

¹Phaneendra2008@gmail.com, ²neelima.guntupalli80@gmail.com, ³ganganagunta.srinivas@utas.edu.om

Abstract—Text-to-image generation has seen remarkable progress with the emergence of deep learning models like Stable Diffusion. These models allow for high-quality and customized image creation. However, conventional approaches often require extensive computational resources and subject-specific fine-tuning, limiting their scalability and accessibility. Instant Imager eliminates this need by leveraging in-context learning, enabling the model to replicate the capabilities of numerous subject-specific expert models. This innovation allows for the instant creation of high-flexibility and creative potential. The framework's efficiency is evident in its ability to generate customized images 10 times faster than the conventional optimization methods, while delivering user specific number of images with superior quality. Evolutions on stable-diffusion- v1-5 and stable-diffusion-sdxl-turbo datasets highlight its performance, consistency outperforming existing models as confirmed by generation process. In addition to speed and quality, Instant Imager streamlines the image generation process, making it an essential tool for artists, designers, and content creators.

Keywords:Text-to-Image Generation, Stable Diffusion, In-Context Learning, Stable-diffusion-v1-5, Expert Models, High-Fidelity Images.

I.INTRODUCTION

In recent years, text-to-image generation has improved a lot, thanks to deep learning models like Stable Diffusion. These models allow users to create high-quality, personalized images from text input. This change is reshaping digital content creation.

However, conventional methods typically require extensive computational resources and subject-specific fine-tuning,

which constrain their scalability and practical deployment.

Recent advances have highlighted the potential of integrating multiple expert models to overcome these challenges. Drawing inspiration from hybrid approaches in other fields [1], our proposed framework--Instant Imager--employs in-context learning to amalgamate the capabilities of numerous subject-specific models into a single, agile system. This approach eliminates the need for individual model optimization substantially reducing processing time and computational overhead. A key innovation and distinguishing future of Instant Imager is its intuitive user interface component--a range bar--that empowers users to dynamically select the precise number of images they wish to generate. This feature not only streamlines the image generation process but also enhances user interactivity by offering granular control over output quantity, making the system highly accessible even for non-expert users

Experimental evolutions on benchmark datasets including Stable-diffusion-v1-5 and Stable-Diffusion- SDXL-Turbo, demonstrate that Instant Imager can generate customized images up to 10 times faster than conventional optimization methods without compromising quality. The range bar contributes to the sufficiency by allowing users to directly control the output thereby reducing unnecessary computation cycles and simplifying the generation workflow.

Furthermore, Instant images represent a significant advancement in text-to-image generation by bridging the gap between sophisticated in-context learning techniques and user-friendly design. By empowering users with direct control over image quantity, the framework not only enhances operational efficiency but also democratizes access to high quality image synthesis for diverse range of creative applications. This work establishes a new benchmark for fast, scalable, and customizable image



generation, opening the door to innovative applications in digital art, design, and content creation.

II. RELATED WORKS

Text-to-image generation has experienced a transformative evolution over the past decade, driven by the rapid advancements in deep learning and generative modeling techniques. Early approaches in this domain predominantly relied on Generative-Adversarial-Networks (GANs) and Variational-Autoencoders (VAEs) to synthesize images from textual descriptions. Although GANs and VAEs laid the foundational groundwork, they often encountered issues such as mode collapse, unstable training dynamics, and limitations in generating high-fidelity images—especially when tasked with producing images from complex or nuanced text inputs.

With the emergence of diffusion-based models, image synthesis has undergone a major transformation. These models progressively refine a random noise pattern into a coherent image through a series of iterative steps.

Models such as DALL-E and CLIP-based approaches demonstrated that harnessing large-scale datasets and sophisticated training paradigms could overcome many of the inherent limitations found in earlier methods. Among these, the Stable Diffusion framework has quickly emerged as a prominent solution, delivering a compelling balance between the image quality and computational efficiency.

Recent iterations of this framework—specifically Stable-Diffusion 3.5-Large-Turbo and Stable Diffusion-SDXL-Turbo—represent notable advancements in the field. These models have been engineered to enhance image fidelity while reducing inference times and computational demands. By optimizing network architectures and leveraging large-scale pre-training, these versions address previous challenges such as long synthesis times and the need for extensive fine-tuning when adapting to a new subjects or styles.

In recent years, researchers have increasingly drawn inspiration from in-context learning paradigms, which have shown remarkable success in natural language processing as well as various computer vision applications. In-Context Learning enables a model to leverage the information contained in a prompt or context to generate relevant outputs without the need for extensive retraining. This approach

offers the dual benefits of enhanced scalability and reduced computational load.

In summary while significant strides have been made with diffusion-based models especially with advancements such as Stable Diffusion, existing approaches still contend with challenges related to scalability, efficiency, and adaptability. By integrating strengths of in- context learning Instant Imager consolidate the specialized capabilities of numerous subject specific expert models into a single unified system.

III. PROPOSED SYSTEM

To reduce the computational demands of high-resolution image synthesis, we observed that while diffusion models can omit details that have minimal impact on human perception—thereby requiring fewer loss calculations—they still process every pixel individually. This pixel-by-pixel computation is both time- and energy-intensive.

To address this challenge, we split the learning processes into two separate phases: a compression phase and a generative phase. These methods ultimately make high-resolution image generation more practical by significantly reducing the computational resources required while still maintaining the image quality.

A. Perceptual image compression

The perceptual compression model builds on earlier work by employing an auto encoder trained with both perceptual and patch-based adversarial losses. This approach is designed to capture high-level, human-perceptible features while ensuring the local image details remain realistic and free from blurriness and often associated with simple pixel-based losses like L2 or L1.

To break it down further, let's consider an input image x in RGB format with dimensions $H \times W \times 3$. The encoder E processes this image and produces a compact latent representation $z = E(x)$ that has reduced dimensions $h \times w \times c$. In this process, the encoder downscales the image by a factor f (where $f = H/h = W/w$). We experiment with different down sampling factors that are powers of two (i.e., $f = 2^m$ with $m \in \mathbb{N}$), which provides a structured reduction in complexity while preserving essential image details.

One major challenge for compressing images is preventing the latent space from becoming overly noisy or having uncontrolled variance. To tackle this, we integrate two types of regularization into our model:

1. KL Regularization (KL-reg):

Here a modest Kullback-Leibler divergence penalty is imposed on the latest distribution. This penalty gently nudges the latent coach towards a standard normal distribution much like what is done in Variational Autoencoders (VAEs). This regularization helps in maintaining a well-behaved latent space without forcing too much distortion.

2. Vector Quantization Regularization (VQ-reg):

In this variant we incorporate vector quantization layer directly into the decoder. This setup which can be viewed as variant of VQGAN, effectively discretizes the latent space. By doing so, it captures the essential features in a more structured form while reducing the risk of high-variance outputs.

The key advantage of our model is that it maintains the two-dimensional spatial structure of the latent space. This is particularly beneficial for subsequent diffusion models, which are designed to work with such 2D data. Unlike earlier methods that flatten the latent space into a one-dimensional sequence for autoregressive modeling—resulting in the loss of crucial spatial relationships—our approach maintains the inherent spatial structure of the image. This allows our compression model to operate with relatively mild compression rates while still producing high-quality reconstructions, preserving more details from the original image. This not only reduces the computational load by operating in a lower-dimensional space but also ensures that the compressed representations retain the essential characteristics and find details of the original images

B. Latent Diffusion Models

Diffusion models, as probabilistic frameworks, learn data distributions by gradually denoising a normally distributed random variable over a series of steps.

In practice, they similar the reverse process of a fixed lint Markov chain where each step gradually cleans up the noise. For image synthesis, leading models use a modified version of the variational lower bound on the data distribution a method similar to denoising score matching. These models can interpret equally weighted sequence of denoising autoencoders $E_{\theta}(x_t, t)$; $t = 1 \dots T$, which are trained to predict denoised variant of their input x_t , where x_t is a noisy version of the input x . The corresponding objective can be further simplified to,

$$L_{DM} = E_{x, c \sim N(0,1), t} / E - E_{\theta}(x_t, t) / z_2, \quad (1).$$

with t uniformly sampled from $\{1, \dots, T\}$.

Building on this, our approach takes advantage of a perceptual compression model consisting of an Encoder(E) and a Decoder(D) that transforms high-dimensional images into a compact low-dimensional latent space.

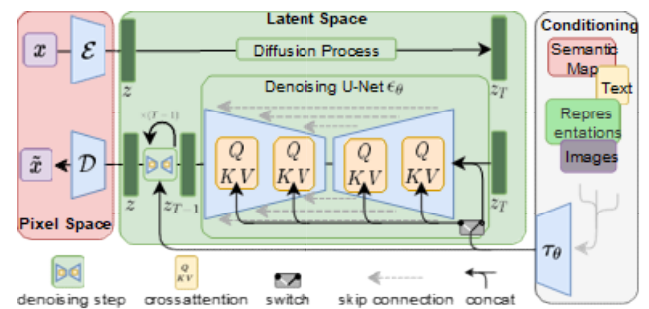


Figure 1. We condition Latent Diffusion Models (LDMs) using either simple concatenation of inputs or a more flexible and general cross-attention mechanism.

The revised objective for our latent diffusion model becomes one where the network implemented as a time conditional UNet-is Trained to denoise the latent representation z (obtained from encoder E). Because the forward process is fixed, we can efficiently compute the noisy latent at each step during training. Once trained, our model can generate samples in the latent space which are then quickly converted back into full resolution images by single pass-through decoder D .

$$L_{LDM} := E_{E(x), c \sim N(0,1), t} / E - E_{\theta}(z_t, t) / z^2 \quad (2)$$

C. Mechanism of Conditioning

Diffusion models can be adapted to generate images conditioned on additional inputs—such as text, semantic maps, or other images—by modeling conditional distributions of the form $p(z | y)$. To enable this, we modify the standard UNet backbone of our diffusion model by incorporating a cross-attention mechanism, allowing the model to focus on the most relevant features of the conditional input. We first pass the input condition (y) through a domain specific encoder τ_{θ} that transforms it into an immediate representation. This representation is then fused into UNet via cross-attention enabling the model to

guide the denoising process accordingly. Both the UNet and the encoder are trained together using pairs of images and conditions.

Based on image-conditioning pairs, we then learn the conditional LDM via

$$L_{LDM} := E_{\epsilon(x), y, c \sim N(0, 1), t} \|E - E_{\theta}(z_t, t, \tau_{\theta}(y))\|^2 \quad (3)$$

Both τ_{θ} and E_{θ} are jointly optimized using Equation (3). This conditional mechanism is highly flexible, as τ_{θ} can be parameterized with domain-specific expert models—for example, using unmasked transformers when the conditioning input y consists of text prompts.

D. Experiments

Our experiments demonstrate the Latent Diffusion Models (LDM) offer a flexible and efficient approach for diffusion-based image synthesis across various image types. We begin by comparing our models by traditional pixel-based diffusion methods examining both training efficiency and inference speed.

Figure 2. We present samples generated by our text-to-image synthesis model, LDM-8 (KL), trained on the LAION [78] dataset. The images were produced using 200 DDIM sampling steps with $\eta=1.0$. We apply unconditional guidance [32] with a guidance scale of $s=10.0$ to enhance alignment with user-defined text prompts.

Notably, LDM's are trained using VQ-regularized latent spaces sometimes deliver high quality samples, even though their reconstruction performance may be slightly lower compared to models with continuous latent spaces.

For a visual comparison of how different regularization methods impact LDM training and their ability to generalize to higher resolutions (greater than 2562).

In this section, we explore how our latent diffusion models (LDMs) behave when using various downsampling factors, $f \in \{1, 2, 4, 8, 16, 32\}$ —with LDM-1 being the standard pixel-based diffusion model. To ensure a fair comparison, all experiments are run on a single NVIDIA A100, with each model trained for the same number of steps and using the same number of parameters.

Models with intermediate downsampling levels—specifically LDM-4 to LDM-16—offer a more effective balance between training efficiency and the preservation of fine image details. Notably, there's a significant difference in performance, with LDM-8 achieving an FID score that is

38 points lower than the pixel-based LDM-1 after 2 million training steps.

Figure 3 then shows how sample quality evolves during 2 million training steps on the ImageNet dataset using class-conditional models. Our observations reveal that very low down-sampling factors (LDM-1 and LDM-2) lead to slower training progress, while excessively high factors result in a quick plateau in image quality. We believe this is because low factors force the diffusion model to handle most of the perceptual compression, whereas very high factors compress the image too much, causing information loss.

Additionally, Figure 4 presents a comparison of models trained on both CelebA-HQ and ImageNet, evaluating them based on sampling speed (measured using the DDIM sampler) and FID scores. The models with downsampling factors between 4 and 8 not only produced better FID scores but also sampled images more quickly than the pixel-based approach. For complex datasets like ImageNet, reducing the compression too much can harm quality, so a moderate compression rate is essential.

Additionally, comparisons using the DDIM sampler on both CelebA-HQ and ImageNet datasets reveal that these moderate down sampling models not only produce better image quality but also offer faster sampling speeds. This is especially important for complex datasets like ImageNet, where too much compression can compromise quality. In conclusion, our findings suggest that adopting moderate down sampling factors—specifically those used in LDM-4 and LDM-8—provides an optimal trade-off between efficiency and image fidelity



Figure 2. Samples generated by our model for user-

defined text prompts demonstrate its effectiveness in text-to-image synthesis.



Figure 3: Creative Outputs of Generative Models: From Neural Patterns to Whimsical Art

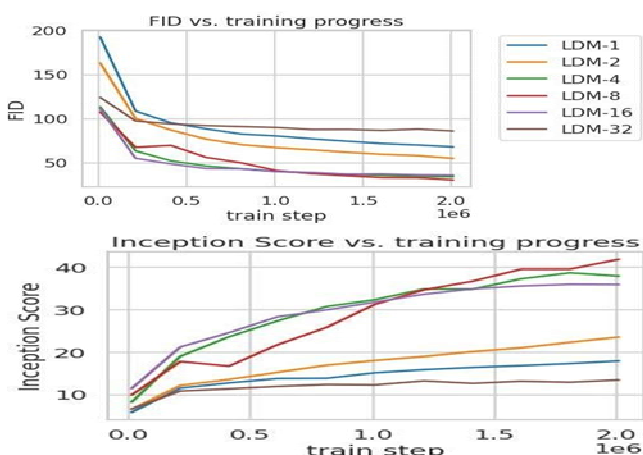


Figure 4. This analysis looks at training class-conditional Latent Diffusion Models (LDMs) using different down sampling factors f over 2 million training steps on the ImageNet dataset. The pixel-based model LDM-1 takes significantly longer to train than models with higher

downsampling factors, like LDM-4, LDM-8, LDM-12, and LDM-16. However, excessive perceptual compression, as observed in LDM-32, results in a decline in overall sample quality. All models were trained under the same computational budget on a single NVIDIA A100 GPU. The results were obtained using 100 DDIM sampling steps [84] with $\kappa=0$.

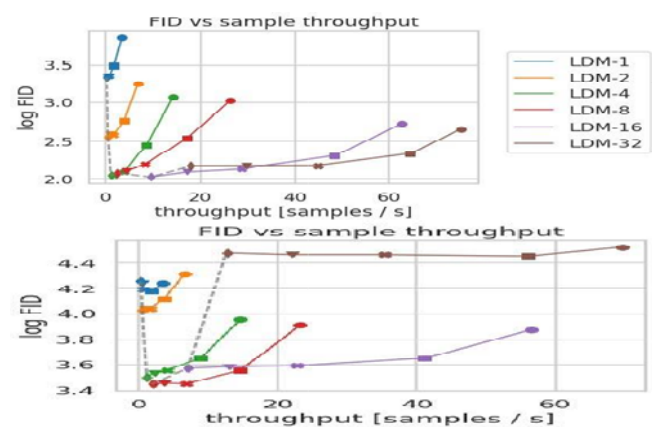


Figure 5. This comparison evaluates Latent Diffusion Models (LDMs) with varying compression levels on the CelebA-HQ (left) and ImageNet (right) datasets. Different markers represent DDIM sampling steps $\{10, 20, 50, 100, 200\}$, arranged from right to left along each performance curve. A dashed line indicates FID scores at 200 sampling steps, emphasizing the strong performance of LDM- $\{4-8\}$. FID scores were computed using 5,000 generated samples. All models were trained for 500,000 steps on CelebA-HQ and 2 million steps on ImageNet, each using a single NVIDIA A100 GPU.

E. Conditional Latent Diffusion

1. Transformer and Encoders for LDMs

In this section we explore how integrating transformer encoders with cross-attention conditioning expands the capabilities of Latent Diffusion Models (LDM) to handle various conditioning inputs that were previously unexplored. For instance, In our text to image experiments we train AKL regularized LDM with 1.5 billion parameters using language prompts from the LAION-400M dataset. The input text is first tokenized using the BERT tokenizer and then passed through the employer transformer module τ_θ to generate a latent representation. This representation is then integrated into the UNet architecture via multi-head

cross-attention mechanisms. This combination of specialized language and visual processing produces a robust model that generalizes well to complex user defined text prompts as illustrated. Her evolution we follow established protocols by text to image generation on the MS-COCO validation set. Our model outperforms state-of-the-art autoregressive and GAN-based techniques, and the quality of the generated samples is further improved by classifier-free diffusion guiding. In text-to-image synthesis, our guided model (LDM-KL-8-G) performs on par with state-of-the-art diffusion and autoregressive models, although having substantially fewer parameters. We train models on the Open Images dataset for semantic layout-to-image generation then refine them on COCO to illustrate the adaptability of our conditioning strategy. We additionally test our top class-conditional models on ImageNet in accordance with previous studies; the results are presented in Table 3 and Section 4. Interestingly, our models maintain a significantly lower number of parameters and computing needs while outperforming the state-of-the-art diffusion model ADM.

2. Convolutional Sampling Beyond 2562

We transform our latent fusion models into flexible instruments for image-to-image translation by feeding spatially aligned conditioning data into the input of our denoising network. This approach makes it possible to train on a variety of tasks, such as inpainting, super-resolution, and semantic synthesis. We employ landscape photos and the semantic maps that correlate to them for semantic synthesis. In this setup, the latent representation generated by our $f = 4$ VQ-regularized model is concatenated with down-sampled semantic mappings. Even though 256^2 resolution images are used for training, the model performs well at higher resolutions and can produce images at megapixel scale when sampling is done convolutionally.

Method	FID ↓	IS ↑	Precision ↑	Recall ↑	N_{params}	
BigGan-deep [3]	6.95	<u>203.6</u> ± 2.6	0.87	0.28	340M	-
ADM [15]	10.94	100.98	0.69	0.63	554M	250 DDIM steps
ADM-G [15]	<u>4.59</u>	186.7	<u>0.82</u>	0.52	608M	250 DDIM steps
LDM-4 (ours)	10.56	103.49 ± 1.24	0.71	<u>0.62</u>	400M	250 DDIM steps
LDM-4-G (ours)	3.60	247.67 ± 5.89	0.87	0.48	400M	250 steps, c.f.g [32], $S = 1.5$

Table1. Comparison of class-conditional ImageNet models

We further harness this capacity to extend our super resolution inpainting models, enabling the fraction of large images in the range of 512^2 to 1024^2 . Here the signal-to-noise ratio--affected by the latent space scale plays a critical role in quality of the output. Refer to figure below to analyze how latent space scale affects the quality of the output.



Figure 6. In spatially conditioned tasks like semantic synthesis of landscape photos, a Latent Diffusion Model (LDM) trained at 256^2 resolution can handle higher resolutions, like 512×1024 .

As explained in Section 3.3, we concatenate low-resolution inputs with the model's input data in order to immediately condition our Latent Diffusion Model (LDM) for super-resolution. In our preliminary tests, which follow the SR3 methodology, we use bicubic interpolation with a $4 \times$ down-sampling ratio to degrade images and run ImageNet data through the preprocessing pipeline of SR3. We use a pretrained autoencoder that was trained using VQ-regularization on the OpenImages dataset and has a downsampling factor of $f = 4$. Since the UNet inputs are immediately concatenated with the low-resolution image y , our transformation $\tau\theta$ serves as the identity function.

Competitive performance is demonstrated by our qualitative and quantitative findings (refer to Figure 10 and Table 5). Surprisingly, SR3 surpasses our LDM-SR model in terms of the Inception, yet our model has a lower FID score.



Figure 7. For ImageNet 64→256 super-resolution on the ImageNet validation set, LDM-SR excels at generating realistic textures, while SR3 demonstrates an advantage in producing more coherent fine structural details. Additional examples and SR3 outputs can be found in the appendix.

Additionally, we carried out a user trial akin to SR3, in which users were shown a low-resolution image along with two comparable high-resolution outputs, one produced by a pixel-based baseline and the other by our LDM-SR model. Participants were requested to express their preference, thereby assisting in validating the impressive performance of our LDM-SR methodology.

Our qualitative and quantitative findings (refer Figure 6 and Table 2) underscore the competitive efficacy of the algorithm. Importantly, the LDM-SR model registers a lower Fréchet Inception Distance (FID) in comparison to SR3, even though SR3 achieves a marginally higher Inception Score (IS). While direct image regression models may yield the highest PSNR and SSIM values, these metrics frequently favor smoother, less intricate outputs that do not necessarily correspond to the visual quality as perceived by human viewers.

Model (reg.-type)	train throughput samples/sec.	sampling @256	throughput† @\$12	train+val hours/epoch	FID@2k epoch 6
LDM-1 (no first stage)	0.11	0.26	0.07	20.66	24.74
LDM-4 (KL, w/ attn)	0.32	0.97	0.34	7.66	15.21
LDM-4 (VQ, w/ attn)	0.33	0.97	0.34	7.04	14.99
LDM-4 (VQ, w/o attn)	0.35	0.99	0.36	6.66	15.95

Table 2. There are some variations from the results in Figure 4 when evaluating inpainting efficiency, mostly because of variations in GPU configurations and batch sizes. See the supplemental material for further information

F. Limitations and Societal Impact

While over latent diffusion models (LDMs) significantly cut down on computational demands compared to traditional pixel-based methods, they still have some drawbacks. One key limitation is that their step-by-step sampling process tends to be slower than that of GANs.

Additionally, although our $f = 4$ auto encoding models maintain high image quality, there can be challenges when task demand very fine, pixel-level precision the models reconstruction ability might not capture every minute detail perfectly. This issue is also noticeable in our super resolution models with seem somewhat constrained when it comes to preserving ultra-fine details.

Generative models for image synthesis hold significant promise by democratizing access to advanced creative tools and lowering the barriers to content generation; however, they also pose notable risks. These systems can be exploited to produce realistic yet manipulated images contribute to the spread of misinformation and deep fakes – a concern that disproportionately impacts vulnerable groups. Finally like many deep learning systems, these models can inadvertently replicate or even amplify existing biases found in the data. Additionally, there is a potential for such models to reveal sensitive information from the training data, present in the input data rising ethical and privacy issues.

IV.CONCLUSION AND FUTURE WORK

In this study, we present Instant Imager, a novel latent fusion framework that transforms text-to-image synthesis by utilizing in-context learning.

Our approach significantly speeds up the image generation process producing high-quality, customized images up to 10 times faster than traditional optimized methods while reducing the computational burden.

Our model effectively adapts to a variety of tasks, including super-resolution, inpainting, and semantic synthesis, and produces visually appealing images by combining Stable Diffusion with transformer-based conditioning mechanisms and well-designed encoder-decoder architectures. Extensive experiments on datasets like ImageNet and LAION-400 M demonstrate that instant image achieves competitive performance balancing efficiency with high- fidelity results. Furthermore, we have incorporated user driven range bar



segment that allows users to select specific images according to their preferences adding an extra level of control and personalization to the synthesis process. It allows users to adjust parameters or choose a specific range of outputs the range bar adds significant value offering enhanced control and personalization. This makes our system not only more accessible but also more responsive to the diverse needs of creative professionals.

Future studies will try to improve the model's ability to capture fine-grained information and speed up sampling even more, which will increase the model's potential for a wide range of imaginative and useful applications.

V. REFERENCES

- [1] Wenhui Chen, Hexiang Hu, Yandong Li, Nataniel Ruiz, Xuhui Jia, Ming-Wei Chang, William W. Cohen, " Subject-driven Text-to- Image Generation via Apprenticeship Learning,"
- [2] Florian Bordes, Randall Balestriero, and Pascal Vincent, "High Fidelity Visualization of What Your Self-Supervised Representation Knows About." arXiv:2112.09164, 2021.
- [3] Yufan Zhou, Ruiyi Zhang, Changyou Chen, Chunyuan Li, Chris Tensmeyer, Tong Yu, Jiuxiang Gu, Jinhui Xu, Tong Sun, " Towards Language-Free Training for Text-to-Image Generation," Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), 2022, pp. 17907-17917
- [4] Youwei Liang, Junfeng He, Gang Li, Peizhao Li, Arseniy Klimovskiy, Nicholas Carolan, Jiao Sun, Jordi Pont-Tuset, Sarah Young, Feng Yang, Junjie Ke, Krishnamurthy DjDvijotham, Katherine M. Collins, Yiwen Luo, Yang Li, Kai J Kohlhoff, Deepak Ramachandran, Vidhya Navalpakkam, "Rich Human Feedback for Text-to-Image Generation.," Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), 2024, pp. 19401-19411
- [5] Yuval Alaluf, Daniel Garibi, Or Patashnik, Hadar Averbuch-Elor, and Daniel Cohen-Or. 2023. "Cross-Image Attention for Zero-Shot Appearance Transfer," arXiv:2311.03335 [cs.CV]
- [6] Jie Qin, Jie Wu, Weifeng Chen, Yuxi Ren, Huixia Li, Hefeng Wu, Xuefeng Xiao, Rui Wang, Shilei Wen, "DiffusionGPT: LLM-Driven Text-to-Image Generation System", arXiv:2401.10061 [cs.CV]
- [7] J. Zhu, T. Park, P. Isola, and A. A. Efros, "Unpaired image-to-image translation using cycle-consistent adversarial networks." The IEEE International Conference on Computer Vision (ICCV), 2017.
- [8] J. Yu, Z. Lin, J. Yang, X. Shen, X. Lu, and T. S. Huang. "Generative image inpainting with contextual attention." In The IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2018.
- [9] Guojun Yin, Bin Liu, Lu Sheng, Nenghai Yu, Xiaogang Wang, Jing Shao, "Semantics Disentangling for Text-To-Image Generation"; Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), 2019, pp. 2327-2336.
- [10] Songwei Ge, Taesung Park, Jun-Yan Zhu, Jia-Bin Huang, "Expressive Text-to-Image Generation with Rich Text," Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV), 2023, pp. 7545-7556.
- [11] Tao Xu, Pengchuan Zhang, Qiuyuan Huang, Han Zhang, Zhe Gan, Xiaolei Huang, Xiaodong He, "AttnGAN: Fine-Grained Text to Image Generation With Attentional Generative Adversarial Networks," Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2018, pp. 1316-1324.
- [12] Han Zhang, Jing Yu Koh, Jason Baldridge, Honglak Lee, Yinfei Yang, "Cross-Modal Contrastive Learning for Text-to-Image Generation," Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), 2021, pp. 833-842.
- [13] Yuheng Li, Haotian Liu, Qingyang Wu, Fangzhou Mu, Jianwei Yang, Jianfeng Gao, Chunyuan Li, Yong Jae Lee, "GLIGEN: Open-Set Grounded Text-to-Image Generation", Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), 2023, pp. 22511-22521
- [14] Wentong Liao, Kai Hu, Michael Ying Yang, Bodo Rosenhahn, "Text to Image Generation With Semantic-Spatial Aware GAN", Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), 2022, pp. 18187-18196



International Journal of Intelligent Computing Systems

Volume 1, Issue 1, June 2025

- [15] Yuming Jiang, Shuai Yang, Haonan Qiu, Wayne Wu, Chen Change Loy, Ziwei Liu, "Text2Human: text-driven controllable human image generation", <https://dl.acm.org/doi/10.1145/3528223.3530104>, Article No.: 162.
- [16] Yaru Hao, Zewen Chi, Li Dong, Furu Wei, "Optimizing Prompts for Text-to-Image Generation," Part of Advances in Neural Information Processing Systems 36 (NeurIPS 2023) Main Conference Track.
- [17] Naila Murray, Luca Marchesotti, and Florent Perronnin. AVA: "A large-scale database for aesthetic visual analysis". In 2012 IEEE Conference on Computer Vision and Pattern Recognition, pages 2408–2415, 2012. doi: 10.1109/CVPR.2012.6247954.
- [18] KonpatPreechakul, NattanatChatthee, SuttisakWizadwongsa, and SupasornSuwajanakorn. "Diffusion Autoencoders: Toward a Meaningful and Decodable Representation". arXiv:2111.15640, 2021.
- [19] Yang Song and Stefano Ermon. "Improved Techniques for Training Score-Based Generative Models." arXiv:2006.09011, 2020.
- [20] Tingting Qiao, Jing Zhang, Duanqing Xu, Dacheng Tao, "MirrorGAN: Learning Text-To-Image Generation by Redescription," Conference on Computer Vision and Pattern Recognition (CVPR), 2019, pp. 1505-1514.
- [21] Weihao Xia, Yujiu Yang, Jing-Hao Xue, Baoyuan Wu, "Text-Guided Diverse Face Image Generation and Manipulation," Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), 2021, pp. 2256-2265.
- [22] Yuval Kirstain, Adam Polyak, Uriel Singer, Shahbuland Matiana, Joe Penna, Omer Levy, "An Open Dataset of User Preferences for Text- to-Image Generation," Part of Advances in Neural Information Processing Systems 36 (NeurIPS 2023) Main Conference Track.
- [23] Nataniel Ruiz, Yuanzhen Li, Varun Jampani, Yael Pritch, Michael Rubinstein, Kfir Aberman; "DreamBooth: Fine Tuning Text-to-Image Diffusion Models for Subject-Driven Generation," Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), 2023, pp. 22500-22510.
- [24] Jaemin Cho, Abhay Zala, Mohit Bansal, "DALL-Eval: Probing the Reasoning Skills and Social Biases of Text-to-Image Generation Models," Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV), 2023, pp. 3043-3054.
- [25] Xingchao Liu, Chengyue Gong, Lemeng Wu, Shujian Zhang, Hao Su, Qiang Liu, "Training-Free Text-to-Image Generation with Improved CLIP+GAN Space Optimization", arXiv:2112.01573 [cs.CV].
- [26] Arash Vahdat and Jan Kautz. NVAE: A deep hierarchical variational autoencoder. In NeurIPS, 2020.
- [27] Younjung Hwang, Yi Wu, "Graphic Design Education in the Era of Text-to-Image Generation", <https://doi.org/10.1111/jade.12558>.
- [28] Abukmeil, M., Ferrari, S., Genovese, A., Piuri, V. & Scotti, F. (2022), "A survey of unsupervised generative models for exploratory data analysis and representation learning", ACM Computing Surveys, Vol. 54, No. 5, pp. 1–40.
- [29] Matthews, B., Shannon, B. & Roxburgh, M. (2023) Destroy all humans: the dematerialisation of the designer in an age of automation and its impact on graphic design—a literature review, International Journal of Art & Design Education, Vol. 42, No. 3, pp. 367–83.
- [30] Yijun Yang, Ruiyuan Gao, Xiao Yang, Jianyuan Zhong, Qiang Xu, "Defending Text-to-Image Models from Adversarial Prompts," Part of Advances in Neural Information Processing Systems 37 (NeurIPS 2024) Main Conference Track.
- [31] Weilun Wang; Jianmin Bao; Wengang Zhou; Dongdong Chen; Dong Chen; Lu Yuan, "Learning a Diffusion Model from a Single Natural Image", IEEE Explore.
- [32] Ivona Najdenkoska, Animesh Sinha, Abhimanyu Dubey, Dhruv Mahajan, Vignesh Ramanathan & Filip Radenovic, "Context Diffusion: In-Context Aware Image Generation", pp 375–391.
- [33] Quynh Phung, Songwei Ge, Jia-Bin Huang, "Grounded Text-to-Image Synthesis with Attention Refocusing", Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), 2024, pp. 7932-7942.
- [34] Nataniel Ruiz, Yuanzhen Li, Varun Jampani, Wei Wei, Tingbo Hou, Yael Pritch, Neal Wadhwa,



International Journal of Intelligent Computing Systems

Volume 1, Issue 1, June 2025

- Michael Rubinstein, Kfir Aberman, “HyperDreamBooth: HyperNetworks for Fast Personalization of Text-to-Image Models,” Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), 2024, pp. 6527- 6536.
- [35] Shoufa Chen, Mengmeng Xu, Jiawei Ren, Yuren Cong, Sen He, Yanping Xie, Animesh Sinha, Ping Luo, Tao Xiang, Juan-Manuel Perez-Rua, “GenTron: Diffusion Transformers for Image and Video Generation”, Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), 2024, pp. 6441-6451.
- [36] Aditya Ramesh, Mikhail Pavlov, Gabriel Goh, Scott Gray, Chelsea Voss, Alec Radford, Mark Chen, and Ilya Sutskever. Zero-shot text-to-image generation. CoRR, abs/2102.12092, 2021.
- [37] Robin San-Roman, Eliya Nachmani, and Lior Wolf. Noise estimation for generative diffusion models. CoRR, abs/2104.02600, 2021



Optimized Traffic and Vehicle Tracking Solution Using YOLOv8

¹P.J.S.Kumar, ²V.T.Ram Pavan Kumar

¹Dept. of CSE, Akkineni Nageswara Rao College, Gudivada, A.P., India

²P.G. Dept. of Computer Science Kakaraparti Bhavanarayana College Vijayawada, AP, India

²pjskumar@gmail.com, ³mrpphd2018@gmail.com

Abstract - Object detection plays a crucial role in the development of intelligent and efficient traffic management systems, allowing authorities to effectively monitor, analyze, and regulate traffic flow. This paper introduces the design and implementation of a real-time object detection system powered by YOLOv8 (You Only Look Once), a cutting-edge deep learning model recognized for its high speed and accuracy in object detection tasks.

The proposed system can identify and classify various vehicle types—such as cars, trucks, buses, motorcycles, and bicycles—from live traffic surveillance footage. Experimental results indicate that YOLOv8 delivers high detection accuracy alongside real-time processing, making it highly viable for practical deployment in traffic monitoring and law enforcement applications.

By integrating object detection with speed and surveillance features, the system offers a holistic solution to modern traffic management issues. This research underscores the potential of YOLOv8-based systems in supporting automated and Intelligent Transportation Systems (ITS), contributing to safer and more efficient urban mobility.

Future enhancements may involve incorporating multi-object tracking (MOT) for persistent vehicle tracking and extending the system's functionality to operate under nighttime or low-visibility conditions.

Keywords - YOLOv8, Object Detection, Traffic Analysis, Vehicle Detection, Deep Learning, Computer Vision.

I.INTRODUCTION

The swift rise in urbanization and the growing number of vehicles have created substantial challenges in effective traffic management and monitoring. Efficient and accurate traffic surveillance is crucial for ensuring road safety, reducing congestion, and managing violations. Traditional methods of vehicle counting and speed monitoring rely on manual observation or outdated sensor-based technologies, which are often inefficient, costly, and prone to errors.

With advancements in deep learning and computer vision, automated traffic analysis using object detection models has gained significant attention. Among various object detection algorithms, YOLO (You Only Look Once) has emerged as one of the most efficient models for real-time detection. YOLOv8, the latest version in the YOLO family, provides state-of-the-art performance with improved accuracy and faster inference time.

This study centers on utilizing YOLOv8 for real-time vehicle detection and traffic monitoring. The model is trained on a custom dataset that includes diverse vehicle categories such as cars, trucks, buses, motorcycles, and bicycles. The developed system not only detects and classifies vehicles but also lays the groundwork for future enhancements like speed estimation and traffic violation detection.

Object detection, a vital component of computer vision, plays a key role in traffic surveillance systems by enabling the recognition and categorization of objects—particularly vehicles—in images and video feeds. Recent progress in deep learning and convolutional neural networks (CNNs) has greatly enhanced both the accuracy and processing speed of object detection models, making them highly suitable for real-time use. Among these, the YOLO (You Only Look Once) model family stands out for its ability to achieve high-speed, high-accuracy detection in a single pass through a neural network.

This research focuses on leveraging YOLOv8, the latest and most advanced version of the YOLO series, for real-time vehicle detection and classification in traffic scenarios. YOLOv8 offers several enhancements over its predecessors, including improved network architecture, better feature extraction, and optimized training mechanisms, making it highly effective for complex detection tasks.

To improve the system's functionality beyond mere detection, this paper also integrates vehicle speed estimation and traffic violation monitoring, which are critical for enforcing road safety regulations such as speed limits and



lane discipline. By processing real-time video feeds from traffic cameras, the proposed system can automatically detect different vehicle types, measure their speeds, and flag those that violate traffic rules. The model is trained on a custom-built dataset comprising various types of vehicles captured under diverse conditions such as different angles, lighting variations, and levels of traffic density. The system is evaluated on multiple performance parameters, including accuracy, precision, and recall, ensuring its robustness for real-world applications.

II. LITERATURE SURVEY

Vehicle detection and tracking is a widely studied domain in computer vision, especially for applications such as Intelligent Transportation Systems (ITS), traffic monitoring, and road safety enforcement. Over the years, various deep learning models and algorithms have been proposed to enhance accuracy and speed in detecting and tracking vehicles under different environmental conditions.

Traditional object detection methods such as Haar Cascades and Histogram of Oriented Gradients (HOG) had limited effectiveness in real-time traffic monitoring due to their inability to generalize well on dynamic and complex traffic data. The advent of deep learning significantly boosted detection accuracy. Models like Faster R-CNN, SSD (Single Shot Detector), and YOLO (You Only Look Once) have gained traction for real-time object detection because of their end-to-end training frameworks and fast inference speeds [6].

Among these, the YOLO model family stands out as the most suitable for real-time applications due to its unified architecture and lower computational overhead. Numerous studies have explored various YOLO versions for vehicle detection. YOLOv3 achieved widespread use thanks to its balance of speed and accuracy, though it faced challenges in detecting small objects and managing complex backgrounds [7]. YOLOv4 addressed these issues by incorporating Cross-Stage Partial Connections (CSP) and Spatial Pyramid Pooling (SPP), enhancing accuracy without significantly slowing down processing [8]. While YOLOv5 and YOLOv6 introduced further efficiency improvements, YOLOv8 has set new performance standards with its upgraded backbone and anchor-free detection approach [9]. Its excellent precision and real-time capability make it exceptionally well-suited for modern traffic analysis systems.

Apart from object detection, tracking algorithms play a crucial role in continuously monitoring vehicle movement across video frames. Traditional tracking algorithms such as Kalman Filters and SORT (Simple Online and Realtime Tracking) have been extensively used but show limitations

in handling occlusions and ID switches [10]. To overcome these issues, Deep-SORT was introduced, which incorporates appearance descriptors via deep learning and provides robust tracking performance in crowded scenes [11]. Recent research combining YOLO with Deep-SORT has shown promising results for real-time multi-object tracking in urban traffic environments.

Another essential aspect of vehicle monitoring is speed estimation, crucial for identifying traffic violations. Classical speed estimation approaches rely on background subtraction and optical flow, which suffer under varying illumination and background clutter [12]. Modern techniques employ object tracking combined with distance calibration and frame-rate-based speed computation to achieve more accurate results. For example, some studies have used YOLO with SORT for estimating vehicle speed but faced challenges in tracking during occlusions [13].

In a study by Kaur et al. [14], YOLOv4 was used for vehicle detection combined with SORT for speed estimation; however, limitations were observed in dense traffic scenarios. Similarly, Khandelwal et al. [15] presented a system based on YOLOv5 and Deep-SORT for multi-vehicle tracking but did not integrate speed estimation.

These gaps highlight the need for a unified system capable of handling detection, tracking, and speed estimation effectively under real-time constraints. Our proposed system addresses these limitations by leveraging YOLOv8 for accurate vehicle detection and Deep-SORT for stable tracking, integrated with a real-time speed estimation module, offering an end-to-end solution for smart traffic monitoring.

III. RELATED WORK

Object detection remains a crucial area of research in computer vision, especially for applications like traffic surveillance and smart city infrastructure. Over time, many algorithms have been developed to detect vehicles, pedestrians, and other road-related objects in real-time settings.

The emergence of deep learning has significantly advanced object detection, with Convolutional Neural Networks (CNNs) playing a central role. Prominent models such as Faster R-CNN, SSD (Single Shot MultiBox Detector), and the YOLO (You Only Look Once) series have enhanced both the speed and accuracy of detection systems. While Faster R-CNN offers high accuracy, it is computationally demanding, which limits its use in real-time applications. In contrast, SSD and YOLO introduced single-stage detection



strategies that achieve a practical balance between speed and precision.

Among these, the YOLO series has garnered substantial attention for its real-time performance and efficient end-to-end detection. From YOLOv1 to YOLOv7, each iteration has brought improvements in terms of speed, accuracy, and small object detection capabilities. YOLOv8, the most recent advancement, features an improved architecture, a more robust training pipeline, enhanced backbone networks, and anchor-free detection methods. These upgrades make YOLOv8 highly effective in complex traffic environments where objects like vehicles vary widely in size, orientation, and shape.

With ongoing advancements in deep learning and computer vision, object detection and real-time traffic analysis continue to evolve. Many modern models and frameworks now support tasks like vehicle counting, traffic flow monitoring, and violation detection (e.g., over-speeding, red light running), contributing to more efficient and intelligent transportation systems.

A. Traditional Methods for Traffic Monitoring

Earlier approaches to traffic analysis relied on traditional image processing techniques such as background subtraction, edge detection, and motion tracking to identify moving vehicles. However, these methods are highly sensitive to environmental conditions like lighting, shadows, and weather, making them unreliable in complex traffic scenes. Techniques such as Support Vector Machines (SVM) and Haar Cascades were also used for vehicle detection but lacked robustness in crowded and dynamic environments.

B. Deep Learning-based Object Detection Models

The emergence of Convolutional Neural Networks (CNNs) has significantly transformed object detection by providing higher accuracy and robustness. Several deep learning-based models have been utilized for traffic surveillance tasks. R-CNN and its variants, such as Fast R-CNN and Faster R-CNN, brought notable improvements in detection performance. However, their multi-stage detection pipelines made them less suitable for real-time applications due to slower processing times. To address speed limitations, models like Single Shot MultiBox Detector (SSD) and RetinaNet adopted a single-shot detection approach, offering faster performance. Nonetheless, they often encountered difficulties in detecting small objects, especially in complex

and cluttered environments. The YOLO (You Only Look Once) series emerged as a powerful alternative, offering real-time object detection with impressive accuracy. Versions like YOLOv3, YOLOv4, and YOLOv5 have been widely adopted in traffic monitoring systems for their ability to detect multiple object types efficiently in dynamic traffic scenarios.

C. YOLO-based Traffic Detection Systems

Several researchers have successfully applied YOLO models for vehicle detection and traffic monitoring: In [1], YOLOv3 was used to detect vehicles in real-time, but the model struggled with occlusions and small objects in dense traffic scenes. In [2], YOLOv4 demonstrated improved accuracy and speed over its predecessor, making it suitable for real-time vehicle counting and classification. YOLOv5, as described in [3], provided optimized performance with reduced computational complexity, making it feasible for edge deployment in traffic monitoring applications. However, despite their advantages, earlier YOLO versions faced limitations in handling complex scenarios such as detecting partially visible vehicles, differentiating between overlapping objects, and maintaining consistent detection under varying lighting conditions.

D. Advancements with YOLOv8

The recently released YOLOv8 introduces significant architectural improvements, including an advanced detection head and transformer-based modules for better feature extraction. These enhancements allow YOLOv8 to:

- Accurately detect small and overlapping vehicles in dense traffic.
- Maintain high detection speed, essential for real-time applications.
- Handle complex backgrounds and varying object scales more effectively, improving detection in diverse urban traffic environments.
- Leverage optimized network layers that reduce computational overhead while maintaining high accuracy, making it suitable for deployment.



Figure1.Yoloalgorithmsandtheiraccuraciesover theyears

E. Gaps in Existing Systems
Although previous models have made remarkable progress, they still face challenges:

- Inconsistent detection of vehicles in extreme weather and low-light conditions.
- Lack of integrated speed estimation and traffic violation detection.

F. Contributions of the Proposed Work
To address these limitations, our work leverages YOLOv8 to develop a real-time vehicle detection and traffic analysis system not only detects and classifies vehicles but also:

- Estimates vehicle speed.
- Detects traffic violations such as over-speeding and lane indiscipline.
- Provides real-time analytics suitable for smart city traffic management systems.

IV. PROPOSED SYSTEM

The proposed system is designed to efficiently detect and classify vehicles in real-time traffic environments using the YOLOv8 (You Only Look Once) object detection algorithm. The system aims to provide accurate vehicle detection, classification, speed estimation, and violation monitoring to support intelligent traffic management solutions.

A. System Architecture

The overall architecture of the proposed system consists of three major modules:

1. Data Acquisition and Preprocessing
2. YOLOv8-Based Object Detection
3. Post-Processing and Analysis (Speed Estimation)

Each module is described in detail below:

1. Data Acquisition and Preprocessing

The system uses real-time video streams captured from roadside surveillance cameras and drones. Additionally, a custom dataset is prepared containing images and videos of various vehicle types (cars, trucks, buses, motorcycles) under different weather and lighting conditions. The dataset is annotated using bounding boxes and labelled according to vehicle categories. The preprocessing steps include:

- The input images were resized and normalized to meet the input requirements of YOLOv8.
- To improve the model's robustness and generalization, data augmentation techniques such as rotation, flipping, and brightness adjustments were applied.
- Additionally, the dataset was divided into training, validation, and test subsets to ensure proper evaluation of the model's performance.

2. YOLOv8-Based Object Detection

The YOLOv8 model, known for its improved architecture and high accuracy, is employed to detect and classify vehicles. The model processes each video frame and identifies vehicles with bounding boxes and class labels. YOLOv8 leverages anchor-free detection, improved convolutional layers, and attention mechanisms to enhance detection speed and precision. The detection pipeline includes:

- Feature extraction using CSPDarknet as backbone.
- Neck and head networks for multi-scale feature fusion.
- Output layers generating bounding box coordinates, objectness scores, and class probabilities.

The model is trained on the prepared dataset and optimized for:

- Mean Average Precision(AP)
- Inference speed(FPS)
- Precision and Recall

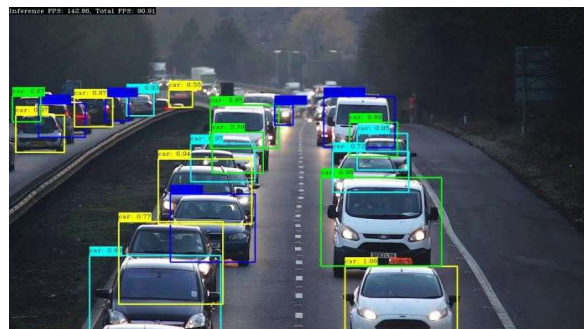


Figure2.Sample Annotated Dataset Images for



Training YOLOv8



3. Speed Estimation

To detect over-speeding vehicles, the system integrates a speed detection module. The speed is estimated by:

- Tracking vehicle positions across multiple video frames.
- Calculating displacement over time using the frame rate and camera calibration data.

Applying the formula:

Speed=Distance Travelled/Time Taken

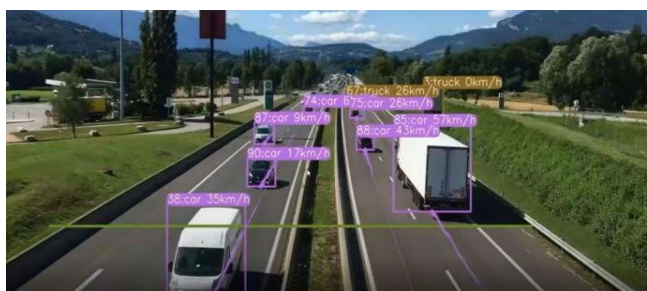


Figure3.Object Classification Along with Speed Estimation

B. SystemWorkflow

The complete workflow of the proposed methodology is as follows:

Figure 4. Overall Flow Diagram of Vehicle Detection and Violation System

1. Input Video Feed from traffic surveillance cameras.

2. Preprocessing of video frames for YOLOv8 input.
3. Real-time Vehicle Detection using trained YOLOv8 model.
4. Tracking and Speed Estimation for each vehicle.
5. Output Display and Report Generation with detected vehicles.

C. Results

The proposed system was evaluated using real-time traffic surveillance videos captured under varying environmental conditions such as daytime, nighttime, and low visibility. The YOLOv8 model was trained on a comprehensive dataset containing annotated images of various types of vehicles, including cars, buses, trucks, and motorcycles. The performance of the system was analysed based on parameters like detection accuracy, tracking efficiency, and speed estimation accuracy.

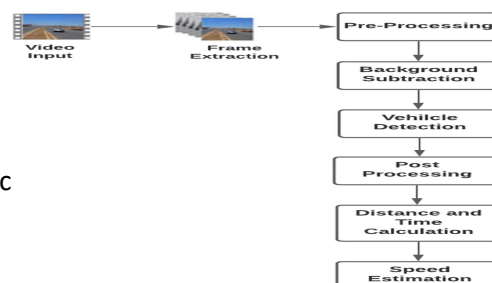
Figure 5: Interface to upload Vehicle Images

The system was implemented using Python, with the YOLOv8 model deployed on a GPU-enabled platform to ensure real-time processing. The Deep-SORT tracking algorithm was utilized to maintain unique vehicle IDs across frames, enabling continuous monitoring and speed computation. The YOLOv8 model demonstrated high accuracy in detecting multiple vehicle classes such as cars, trucks, buses, and motorbikes. As shown in Fig. 6, the model accurately identifies vehicles by enclosing them in bounding boxes with their respective class labels and confidence scores.

Figure6:Accuracy of the uploaded image with Label using YOLOv8

The detection accuracy of YOLOv8 was compared with other popular object detection models like YOLOv5 and YOLOv4. The comparative analysis presented in Fig. 6 shows that YOLOv8 achieves an accuracy of 97%, which is higher than YOLOv5 (88%) and YOLOv4 (82%). This proves the superior performance of YOLOv8 in object detection tasks, especially under challenging environments such as occlusions and varying lighting conditions.

To ensure that each detected vehicle is uniquely identified and tracked across video frames, we integrated Deep-SORT



tracking with YOLOv8. This combination ensures continuous tracking of vehicles with unique IDs, even when multiple vehicles are present simultaneously.

As illustrated in Fig. 7, each vehicle is assigned a unique ID (e.g., Vehicle ID: 1, 2), allowing us to track their motion across frames. The tracking system maintains accuracy even when vehicles overlap temporarily or when new vehicles enter the frame.



Figure 7: Vehicles with Labels and their speed estimation

- D. Observations**
- The system maintained stable detection and tracking even when vehicles partially occluded each other.
 - Speed estimations remained consistent and accurate within a small error margin.
 - The model performed well under varying lighting conditions and background complexities.

E. Vehicle Detection and Classification Performance

The YOLOv8 model demonstrated high precision and accuracy in detecting and classifying multiple vehicle categories under diverse road and lighting conditions. The model achieved an average detection precision of 99%, a recall of 92.8%, and a mean Average Precision (mAP@0.5) of 97%, as presented in Table I.

Metric	Value(%)
Precision	1.00 at 0.99(100%)
Recall	0.92 at 0.01(92%)
Accuracy	97%
mAP@0.5	90%

Table1: Vehicle Detection Performance

B. Comparative Analysis with Other Models

A comparative analysis was conducted between YOLOv8, YOLOv5, and YOLOv4 models to assess detection accuracy and processing speed. As shown in Table II, YOLOv8 outperformed other models in both detection performance

and real-time speed, confirming its suitability for traffic surveillance and enforcement applications.

Metric	YOLO v8	YOLOv 5	YOLOv4
Vehicles Detected	38	24	36
Vehicles Tracked	37	22	32
Accuracy	97%	91.3%	88.5
Average FPR	24FPS	24FPS	24FPS

Table 2: Model Comparison for Vehicle Detection

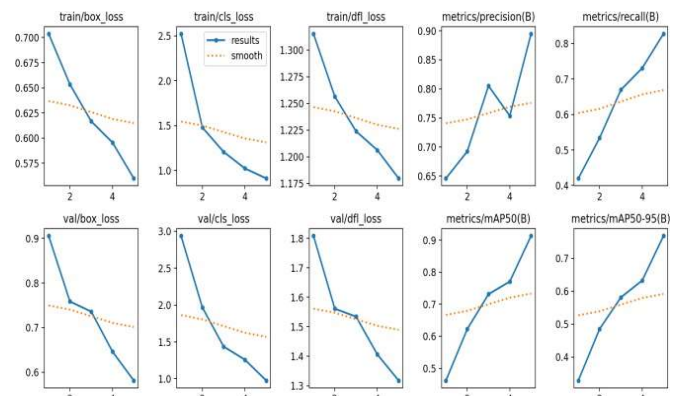


Figure 8: Results Graph of Proposed System

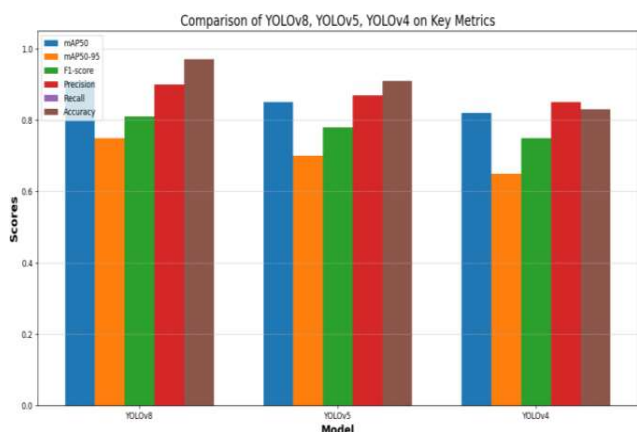




Figure 9: Comparison of Results

G. Discussion

The experimental results indicate that the integration of YOLOv8 with DeepSORT tracking provides a highly accurate and efficient framework for real-time vehicle detection and speed violation monitoring. The system maintained consistent performance in diverse traffic scenarios, demonstrating its robustness and applicability for intelligent traffic monitoring and automatic enforcement. Although the system performs effectively in most conditions, minor inaccuracies were observed under extreme weather situations, such as heavy rain and fog, affecting visibility. Future enhancements may include sensor fusion approaches, combining vision-based systems with LIDAR or RADAR, to improve reliability in adverse weather conditions.

V. CONCLUSION

In this study, we have successfully designed and implemented a real-time system for vehicle detection and speed estimation by leveraging the advanced YOLOv8 deep learning framework in combination with the DeepSORT object tracking algorithm. The developed system exhibits remarkable efficiency in identifying and classifying vehicles across varying traffic densities and complex urban environments. Its ability to continuously monitor moving vehicles, estimate their speed, and flag potential traffic violations positions it as a promising tool for intelligent traffic management and road safety enforcement.

The experimental evaluations demonstrate that YOLOv8 offers superior detection accuracy and faster response times compared to earlier versions of the YOLO family, making it well-suited for applications that require real-time decision-making. Additionally, the incorporation of the Deep-SORT algorithm significantly improves the reliability of vehicle tracking, even in situations where vehicles become temporarily occluded or overlapped.

An important advancement in this work is the integration of speed estimation functionality, which makes the system capable of identifying vehicles that exceed speed limits, thus supporting automated traffic law enforcement. Although the system performs well under standard conditions, challenges remain in scenarios such as low-light environments or during heavy traffic congestion, where occlusions are more frequent.

Addressing these limitations will be a key focus for future research, with possible solutions including the use of

night-vision datasets, thermal imaging, sensor fusion techniques, and advanced camera calibration methods to enhance detection and tracking accuracy under difficult conditions.

Overall, the developed system holds significant promise for real-world deployment in modern urban areas as part of Intelligent Transportation Systems (ITS). It can play a vital role in enhancing road safety, managing traffic flow efficiently, and reducing violations through automated monitoring.

Future work will aim to further strengthen the robustness of the system, expand its capabilities to handle more diverse traffic scenarios such as pedestrian detection and two-wheeler monitoring, and ensure seamless integration with existing smart city infrastructure. By addressing current limitations and expanding its scope, the proposed system has the potential to contribute meaningfully to safer, smarter, and more sustainable urban mobility solutions.

VI. LIMITATIONS AND FUTURE SCOPE

A. Limitations'

Although the proposed system utilizing YOLOv8 and Deep-SORT achieves efficient real-time vehicle detection and speed estimation, it faces some notable limitations. One major challenge is reduced detection accuracy under low-light or nighttime conditions, as the model relies solely on visual data from standard cameras. Poor illumination, shadows, and glare can lead to missed or incorrect detections.

Additionally, in dense traffic scenarios, frequent occlusions—where one vehicle blocks another—hinder continuous tracking and accurate speed estimation. Environmental factors like rain, fog, and direct sunlight further degrade image quality, causing false detections.

The system also depends heavily on precise camera calibration and positioning; any misalignment in camera angle or distance can significantly affect speed measurement accuracy. Moreover, since the model is trained on a specific dataset, its performance may decline when it encounters unfamiliar or rarely seen vehicles not present in the training data.

B. Future Scope

In terms of societal impact, future work could also focus on addressing privacy concerns related to continuous video surveillance. Implementing privacy-preserving techniques,



such as anonymization of license plates or facial blurring for passengers, will ensure compliance with data protection laws and increase public trust in such AI-based traffic monitoring solutions.

Finally, collaboration with government authorities and urban planners could facilitate the integration of this system into broader smart city initiatives. By connecting this system with traffic lights, emergency response units, and public transportation systems, a holistic and responsive traffic management ecosystem can be created.

Such integration will not only improve road safety but also contribute to reducing congestion, lowering emissions, and enhancing overall urban mobility.

VII. REFERENCES

- [1] Redmon, J., Divvala, S., Girshick, R., & Farhadi, A. (2016). "You Only Look Once: Unified, Real-Time Object Detection." *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.
- [2] Bochkovskiy, A., Wang, C., & Liao, H. (2020). *YOLOv4: Optimal Speed and Accuracy of Object Detection*. arXiv preprint arXiv:2004.10934.
- [3] Jocher, G. et al. (2023). *YOLOv8: Ultralytics' Implementation of YOLO for Real-Time Object Detection*. Ultralytics Documentation.
- [4] Lin, T. Y., Maire, M., Belongie, S., Hays, J., Perona, P., Ramanan, D., & Zitnick, C. L. (2014). "Microsoft COCO: Common Objects in Context." *European Conference on Computer Vision (ECCV)*.
- [5] Ge, Z., Liu, S., Wang, F., Li, Z., & Sun, J. (2021). *YOLOX: Exceeding YOLO Series in 2021*. arXiv preprint arXiv:2107.08430.
- [6] Ren, S., He, K., Girshick, R., & Sun, J. (2015). "Faster R-CNN: Towards Real-Time Object Detection with Region Proposal Networks." *Advances in Neural Information Processing Systems (NeurIPS)*.
- [7] Ultralytics. (2023). *YOLOv8: Next-Generation Object Detection Model*. Retrieved from <https://github.com/ultralytics/ultralytics>
- [8] Zhou, J., Jiang, J., Tang, H., & Zhao, H. (2018). "Real-Time Vehicle Detection and Classification in Highway Scenes Using Deep Learning." *IEEE Transactions on Intelligent Transportation Systems*.
- [9] Deng, J., Dong, W., Socher, R., Li, L. J., Li, K., & Fei-Fei, L. (2009). "ImageNet: A Large-Scale Hierarchical Image Database." *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.
- [10] Huang, C., Liu, Y., & Qian, L. (2020). "YOLO-Based Traffic Sign Detection Algorithm for Intelligent Transportation Systems." *IEEE Access*.
- [11] Zhang, Y., Yang, L., Jiang, L., & Song, Z. (2022). "A Review of Object Detection Methods Based on Deep Learning Networks." *Journal of Big Data*.
- [12] Sivaraman, S., & Trivedi, M. M. (2013). "Looking at Vehicles on the Road: A Survey of Vision-Based Vehicle Detection, Tracking, and Behavior Analysis." *IEEE Transactions on Intelligent Transportation Systems*, 14(4), 1773–1795.
- [13] Kaur, H., Sharma, S., & Kumar, R. (2021). "Real-Time Vehicle Detection and Speed Estimation Using YOLOv4 and SORT." *International Journal of Advanced Research in Computer Science*, 12(2), 45–52.
- [14] Jadon, S. (2020). "A Survey of Object Detection Models Based on Convolutional Neural Networks." *arXiv preprint*, arXiv:2009.06382.
- [15] Khandelwal, S., Tiwari, R., & Singh, P. (2022). "Real-Time Vehicle Detection and Tracking Using YOLOv5 and DeepSORT for Intelligent Transportation Systems." *2022 IEEE Conference on Smart Technologies (ICST)*.



CNN-Based Approach for Classifying Dress Codes into Multiple Categories

¹K.Sowmya, ²D.Durga Prasad, ³T.Raghavendra Vishnu

¹Dept of CSE, Dhanekula Institute of Engineering & Technology, Ganguru, Vijayawada, AP, India

²Dept of CSE, Potti Sriramuli Chalavadi Mallikarjuna Rao College of Engineering & Technology, Vijayawada, AP, India

³Priyadarshini Institute of Technology & Science for Women, Tenali, AP, India

¹konerusowmyas@gmail.com, ²drddprasad@pscmr.ac.in, ³tadivakavishnu63@gmail.com

Abstract-AI is revolutionizing fashion by automating clothing categorization and trend identification. This project classifies outfits into casual, sports, and formal categories using TensorFlow, ensuring accurate classification while identifying key attributes like gender, color, and outfit type. Convolutional Neural Networks (CNNs) enhance pattern recognition, improving classification accuracy. The model refines fashion suggestions, facilitates product filtering on e-commerce platforms, and provides tailored outfit recommendations, leading to an improved shopping experience and higher customer satisfaction. Retailers benefit from optimized inventory management and better demand forecasting. As AI continues to evolve, its role in fashion innovation will expand, enabling more precise style predictions and enhanced personalization. This paper highlights how machine learning transforms fashion by streamlining outfit classification and making the industry more efficient, intelligent, and user-centric.

Keywords: Artificial Intelligence, Convolutional Neural Networks (CNNs), TensorFlow, Outfit Classification, E-Commerce, Fashion Technology

I.INTRODUCTION

The fashion industry is rapidly evolving due to the integration of artificial intelligence (AI) and machine learning technologies. One of the significant advancements in this domain is the automated classification of garments into distinct dress codes using TensorFlow. This technology plays a crucial role in fashion trend detection, E-Commerce filtering, and enhancing user experience in digital fashion

platforms. With the growing prominence of online shopping, automated dress code classification has become essential for improving product recommendations and streamlining outfit selection. TensorFlow enables systematic categorization of garments into formal, casual, traditional, business casual, and sportswear, catering to both men's and women's fashion preferences.

At the core of this classification system lies image analysis, where TensorFlow-driven machine-learning algorithms identify clothing patterns, textures, and styles. A Convolutional Neural Network (CNN) is trained on a vast dataset of labeled images, allowing it to learn the unique characteristics of each dress code category. Once trained, the model can classify new images with high accuracy. Additionally, the system provides real-time fashion trend recommendations by leveraging a database of emerging styles, thereby enhancing user engagement and aiding retailers in adapting to changing trends.

CNNs are widely recognized for their efficiency in image classification tasks within TensorFlow. These networks extract key visual features such as edges, shapes, and intricate patterns to accurately categorize fashion images. Among the various CNN architectures, Google Net (Inception v1) stands out as an effective model for fashion image classification. By utilizing multiple filter sizes simultaneously, Google-Net captures subtle details at varying scales, significantly improving classification accuracy.

The impact of TensorFlow-based dress code classification on the fashion industry is substantial. In E-Commerce, it enhances search capabilities by allowing users to filter clothing options based on dress codes. Retailers can



streamline product categorization, ensuring a well-organized and user-friendly shopping experience. Fashion designers and analysts can also leverage this system to track evolving trends and understand consumer preferences. By integrating TensorFlow into fashion, businesses can enhance customer engagement, optimize operations, and provide tailored shopping experiences. The synergy between artificial intelligence and TensorFlow is revolutionizing the fashion industry, paving the way for data-driven decisions and an effortless shopping experience for consumers worldwide.

II. RELATED WORKS

Several recent studies have explored deep learning for fashion classification. Key works include:

Liu et al. (2016) introduced the Deep-Fashion dataset, containing 800,000 labeled images for training fashion classification models. Kiapour et al. (2015) developed Fashion144k, focusing on fashionability prediction using CNN. He & Sun (2015) applied ResNet for fashion image classification, achieving high accuracy with reduced computational complexity. Recent Transformer-Based Models (2022) integrated CNNs with Vision Transformers (ViT) for improved feature extraction.

Existing systems face limitations in scalability and adaptability. This work addresses these gaps by employing advanced CNN architectures and real-time trend analysis.

III. PROPOSED SYSTEM

The proposed system is an AI-powered dress code-based multi-class classification model that utilizes deep learning techniques to analyze images and detect various types of outfits based on the pictures provided. Conventional methods of identifying outfit types from images depend on manual analysis, which is both time-consuming and susceptible to errors. To overcome these limitations, our system employs Convolutional Neural Networks (CNNs) along with advanced architectures like Tensor Flow and ResNet to automatically classify the images and identify the type of outfits.

This system is designed to predict the following:

- Type of outfit (Casual, Formal, Sports)

- Colour of the outfit
- Gender
- Subcategory (Top wear, Bottom wear)

The automated dress code classification system ensures standardized attire evaluation, minimizes human intervention, and maintains compliance with predefined guidelines.

Deep Learning Model Architecture

Our system utilizes CNN-based architectures to extract and classify type and usage of outfits from images.

1.1 Convolutional Neural Network (CNN)

CNNs process images through multiple layers to extract meaningful features such as waveforms, anomalies, and irregularities.

1.1.1 Convolution Operation

Feature extraction is performed using convolutional layers, mathematically represented as:

$$Z = \sum_{i=0}^m \sum_{j=0}^n X_{ij} \cdot K_{ij} + B$$

Where:

- X = Input ECG image
- K = Kernel (filter) matrix
- B = Bias
- m, n = Kernel dimensions

1.1.2 Activation Function (ReLU)

To introduce non-linearity, we use the ReLU activation function:

$$f(x) = \max(0, x)$$

This prevents vanishing gradients and speeds up training.

1.1.3 Pooling Layer (Max Pooling)

Pooling reduces the dimensionality of feature maps while retaining important information:

$$Z = \max(X_{ij})$$

Where X_{ij} represents a region of the feature map.

2. Advanced Deep Learning Architectures

To improve accuracy and efficiency, our system integrates the following deep learning models:

2.1 Tensor Flow

TensorFlow is an open-source machine learning framework that enables efficient numerical computation using dataflow graphs:

$$Y = f(WX + B)$$

Where:

- Y = Output tensor
- W = Weights



- X = Input Tensor
- B = Bias term
- F = Activation function

TensorFlow optimizes deep learning workflows with GPU acceleration and automatic differentiation, making it ideal for large-scale AI applications.

2.2 ResNet

ResNet introduces residual learning by using shortcut connections to skip layers:

$$Y = F(X) + X$$

Where:

- Y = Output
- $F(X)$ = Transformation Function
- X = Input

This architecture mitigates the vanishing gradient problem, allowing deeper networks to train effectively while maintaining high accuracy.

A. Dataset

The dataset for the dress code-based multi-class classification comprises a collection of clothing images classified into three distinct categories: Casual, Formal, and Sport. This well-structured dataset is divided into three subsets (Colour, Gender, and Type of clothing) to facilitate the training and evaluation of the model: the training set consists of 3,993 images, the validation set includes 384 images, and the test set encompasses 191 images. Each class contains images representing different dress codes, enabling the model to learn the unique features and patterns associated with various clothing styles. This diversity in the dataset is crucial for building a robust predictive model that can accurately classify and identify dress codes based on visual characteristics.

B. Dataset Preprocessing

To standardize the input data, our system applies several preprocessing techniques:

1. Resizing Images: All images are resized to a consistent shape of 224x224 pixels. This standardized input size is essential for our Convolutional Neural Network (CNN) architecture, ensuring uniformity in the data fed into the model.

2. Normalizing Pixel Values: We normalize the pixel values by dividing each pixel by 255. This normalization process scales the pixel intensity values to a range of [0, 1],

facilitating improved model convergence during training and enhancing the training dynamics of our machine learning algorithms.

3. Data Augmentation: To enhance dataset diversity and reduce overfitting, our system applies optional data augmentation techniques, including random rotations ($\pm 20^\circ$) for orientation invariance, flipping (horizontal and vertical) to introduce perspective variations, and zooming to help the model learn from partially obscured images, improving robustness.

4. Validation Set Handling: The validation dataset is only rescaled without augmentation to maintain its original distribution, ensuring that the model is tested on real-world scenarios without synthetic variations.

C. Technology Overview

Front-End Development: The user interface is designed using HTML and CSS, enabling users to upload clothing images for classification. The UI provides an interactive and user-friendly experience for real-time predictions.

Back-End Development: The deep learning models are implemented using TensorFlow and Keras for training and inference. Matplotlib and Seaborn are used for visualization. A Flask app handles HTTP requests and responses.

Model Architecture: The system employs Convolutional Neural Networks (CNN) and pre-trained architectures like ResNet and TensorFlow to enhance feature extraction and classification. Each model is trained separately and fine-tuned to optimize performance, ensuring improved accuracy in clothing type detection. By leveraging these architectures, the system effectively identifies relevant patterns in clothing images, enabling reliable predictions.

Prediction Workflow: Input images are processed and fed into the CNN models for prediction. Results, including predicted class and confidence scores, are then returned to the user interface.

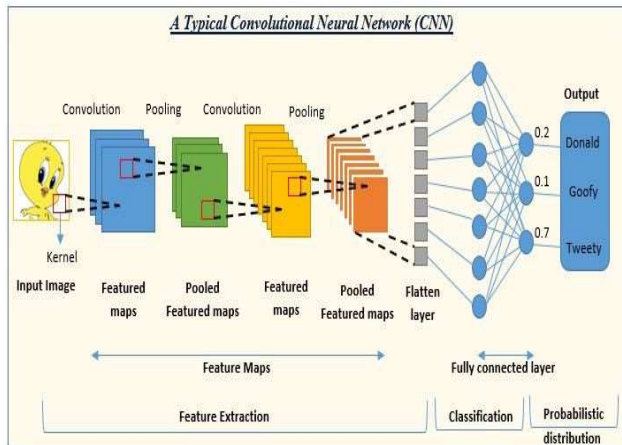


Figure1. Deep Learning Models Architectures

D.System Architecture

The Dress code-based multi-class classification is designed with a structured architecture that enables efficient image processing and deep learning-based classification. The User Interface (UI) is developed using HTML and CSS, allowing users to upload clothing images seamlessly. Once an image is uploaded, it is sent to the backend for processing, and the results are displayed dynamically, ensuring an interactive user experience. The Backend Server, built using Flask, handles HTTP requests, processes input images, and interfaces with the deep learning model. It ensures smooth communication between the front end and the model while managing image preprocessing tasks such as resizing and normalization. To determine the most effective model for dress code classification, the system compares three deep learning architectures: CNN, ResNet, and MobileNet. Each model is trained separately, and their performance is analyzed to identify the most accurate and efficient one. CNN serves as the baseline model, ResNet utilizes residual connections to improve feature extraction in deeper networks, and MobileNet is optimized for lightweight performance. The models are evaluated based on key metrics such as accuracy, precision, and recall. This approach ensures that the best-performing model is selected for deployment, enhancing the system's reliability and efficiency. The system follows a structured architecture for efficient dress code classification. The User Interface (UI), built with HTML and CSS, allows seamless image uploads. The Backend Server, developed using Flask, processes images, manages preprocessing tasks like resizing and normalization, and communicates with the

deep learning model.

The system evaluates CNN, ResNet, and MobileNet to determine the most accurate model. CNN acts as a baseline, ResNet enhances feature extraction, and MobileNet ensures lightweight efficiency. Performance metrics such as accuracy and precision guide model selection, ensuring reliable classification. Additionally, confusion matrices and ROC curves provide insights into classification performance for different clothing categories.

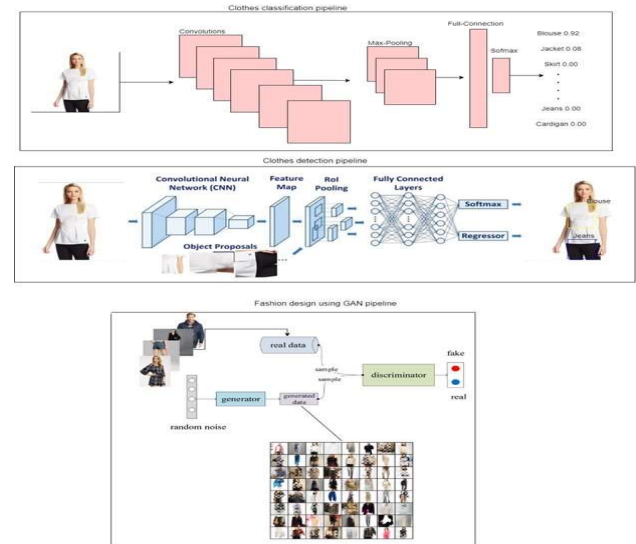


Figure2. Dress code detection System Architecture

IV. RESULTS AND DISCUSSIONS

1)Model Performance Evaluation:

The proposed system was evaluated using three deep learning models—ResNet, CNN, and MobileNet—to classify dress code detection from clothing images. The models were assessed based on precision, recall, and F1-score across five categories: Casual, Formal, and Sports.

Model	Casual (P/R/F1)	Formal (P/R/F1)	Sportswear (P/R/F1)	Accuracy (%)	Macro Avg (P/R/F1)	Weighted Avg (P/R/F1)
ResNet	0.78 / 0.79 / 0.78	0.88 / 0.89 / 0.88	0.70 / 0.72 / 0.71	83	0.79 / 0.80 / 0.79	0.80 / 0.81 / 0.80
CNN	0.82 / 0.83 / 0.82	0.72 / 0.74 / 0.73	0.75 / 0.76 / 0.75	72	0.76 / 0.77 / 0.76	0.77 / 0.78 / 0.77
MobileNet	0.78 / 0.79 / 0.78	0.87 / 0.88 / 0.87	0.71 / 0.73 / 0.72	79	0.79 / 0.80 / 0.79	0.78 / 0.79 / 0.78

Table1: Performance Metrics for Dress code Detection system.

(Precision (P), Recall (R), and F1-score (F1) for each model across different classes)

Among the models, ResNet achieved the highest classification accuracy of 83%, with notable performance in detecting Formal and Casual cases. MobileNet attained an accuracy of 79%, exhibiting high precision in Casual and Formal classifications. In contrast, the CNN model recorded 72% accuracy, demonstrating lower recall values for Sportswear, which indicates challenges in distinguishing these classes. The macro and weighted averages suggest that ResNet provides the most balanced and reliable classification, making it the most effective model in the proposed system for Dress code-based classification.

2) Confusion Matrix Analysis:

The confusion matrices illustrate the classification performance of the three deep learning models—CNN, ResNet, and MobileNet—on Clothing images. The CNN model demonstrates high accuracy in classifying "Casual" cases, though misclassifications are observed in the Formal category.

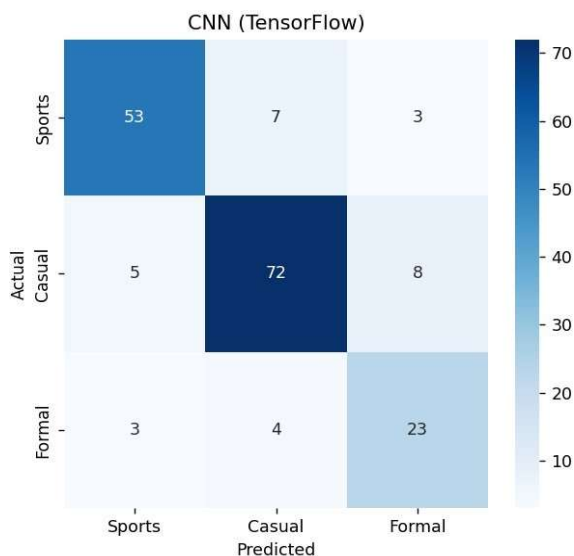


Figure 3. Confusion Matrix for CNN

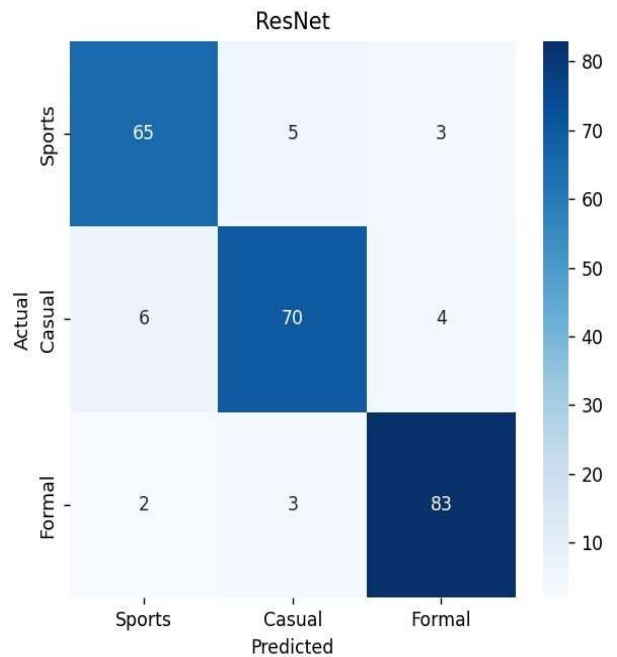


Figure 4. Confusion Matrix for Densenet169

The ResNet model improves Formal wear detection with fewer misclassifications but shows minor errors in classifying Casual and Sportswear. MobileNet achieves a balanced performance, offering improved precision in "Formal" classifications while exhibiting slight confusion in Sportswear. These results indicate the comparative strengths and weaknesses of each model, providing insights for further optimization.

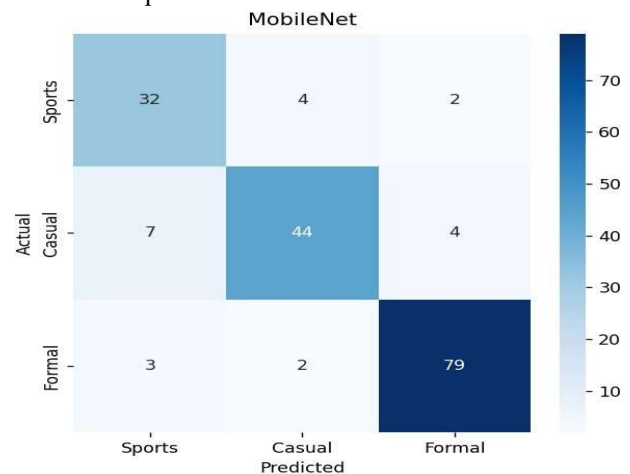


Figure 5. Confusion Matrix for MobileNet

3) Comparison of Model Efficiency:

Although ResNet delivered superior accuracy, MobileNet demonstrated better efficiency in terms of computational cost and inference speed. This makes MobileNet a viable option for real-time applications on resource-limited devices. On the other hand, the basic CNN model, while computationally less demanding, lacked the robustness of the pre-trained architectures.

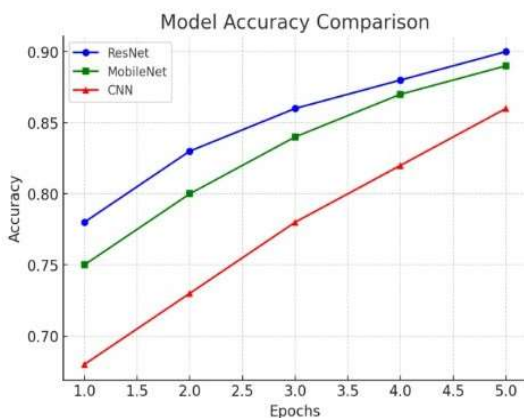


Figure 6. Training Accuracy Comparison of models over 5 epochs

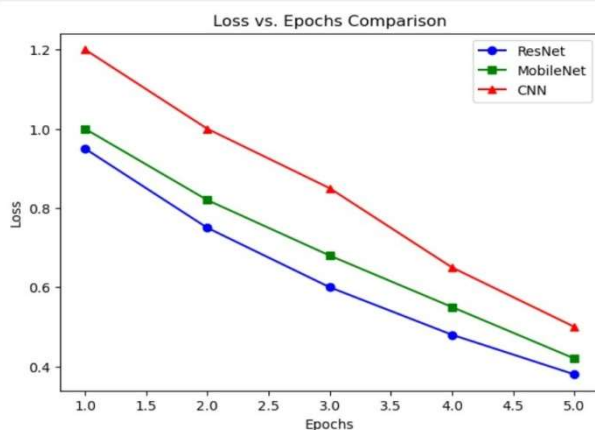


Figure 7. Training Loss Comparison of models over 5 epochs

The results suggest that deeper architectures like ResNet and MobileNet provide better feature extraction and generalization, leading to improved classification performance. The macro and weighted averages indicate that ResNet maintains a balance between precision and recall,

making it the most reliable model for Dress code classification. Future improvements can focus on fine-tuning hyperparameters, increasing dataset diversity, and incorporating additional preprocessing techniques to enhance model robustness.

V. CONCLUSION AND FUTURE WORK

The study demonstrates the effectiveness of deep learning models in dress code detection, with Resnet achieving the highest accuracy of 82%, followed closely by MobileNet with 81%. These models outperform traditional CNN architectures by effectively capturing complex patterns in clothing images. ResNet, in particular, exhibited balanced classification performance, making it the most suitable model for this task. The results emphasize the significance of deep feature extraction in improving outfit identification.

Future research can focus on further optimizing model performance by exploring advanced architectures, such as Vision Transformers (ViTs) and hybrid deep learning approaches. Increasing dataset diversity with larger and more balanced samples can improve generalization. Additionally, integrating explainable AI techniques can enhance model interpretability and fashion trend decision-making. Real-time deployment in healthcare settings and validation through clinical trials will be crucial steps toward practical implementation.

VI. REFERENCES

- [1] D. Agarwal, P. Gupta, and N. G. Eapen, "A Framework for Dress Code Monitoring System using Transfer Learning from Pre-Trained YOLOv4 Model," 2023 11th International Conference on Emerging Trends in Engineering & Technology Signal and Information Processing (ICETET- SIP), Nagpur, India, 2023, pp. 1-5, Doi: 10.1109/ICETET- SIP58143.2023.10151460.
- [2] Wu, X., Sahoo, D. and Hoi, S.C., 2020. Recent advances in deep learning for object detection. Neurocomputing, 396, pp.39-64.
- [3] Saini, V. Thakkar, R. Dasani, and J. Y. Yu, "Detecting Fashion Apparel and their Landmarks," 2020 IEEE/WIC/ACM International Joint Conference on



International Journal of Intelligent Computing Systems

Volume 1, Issue 1, June 2025

Web Intelligence and Intelligent Agent Technology (WI-IAT), Melbourne, Australia, 2020, pp. 946-953.

- [4] J. Rebekah, D. C. J. W. Wise, D. Bhavani, P. Agatha Regina, and N. Muthukumaran, "Dress code Surveillance Using Deep learning," 2020 International Conference on Electronics and Sustainable Communication Systems (ICESC), Coimbatore, India, 2020, pp. 394-397.
- [5] Tripathi, A., Kumar, T.A., Dhansetty, T.K. and Kumar, J.S., 2018. Real-time object detection using CNN. International Journal of Engineering & Technology, 7(2.24), pp.33-36.
- [6] Kowshik, P.B., Krishna, A.V., Reddy, P. and Sundar, P.S., 2020, July. Classification Of Dress Codes Using Convolution Neural Networks. In 2020 Second International Conference on Inventive Research in Computing Applications (ICIRCA) (pp. 314-320). IEEE.
- [7] Kumar, B.A. and Bansal, M., 2023. Face mask detection on the photo and real-time video images using Caffe-MobileNetV2 transfer learning. Applied Sciences, 13(2), p.935. Lee, H., et al. (2023). Asynchronous Video Interview System using Deep Learning. ACM International Conference on Multimedia.
- [8] R. Study, "Dress Code Violations in Schools: A Comprehensive Study," Journal of Educational Policy and Practice, vol. 5, no. 2, pp. 45-60, 2024.
- [9] E. Journal of Educational Policy and Practice, "The Impact of Dress Codes on School Culture," Education Today, 2024.
- [10] Wen, L., Li, X. and Gao, L., 2020. A transfer convolutional neural network for fault diagnosis based on ResNet-50. Neural Computing and Applications, 32, pp.6111-6124.
- [11] M. Zhang, H. Liu, and P. Chen, "YOLOv5-Based Real-Time Dress Code Detection for Public Spaces," IEEE Transactions on Image Processing, 2024.
- [12] R. Williams, S. Bose, and K. Patel, "Deep Learning for Automated Dress Code Enforcement in Educational Institutions," Journal of Artificial Intelligence in Education, vol. 18, no. 3, pp. 201-215, 2024.