# FuelGuard: Fuel Consumption Anomaly Detection and Visual Verification in Logistics Using Isolation Forest, CBIR, and OCR

**Sigit Auliana\*[1], Basuki Rakhim Setya Permana[2], Mochammad Darip[3], Sujan Chandra Roy[4]**

[1,2,3]Computer Science, Universitas Bina Bangsa, Indonesia
[4]Institute of Information and Communication Technology (IICT), Chittagong University of
Engineering and Technology, Bangladesh

Email: [1]pasigit@gmail.com

## Abstract

Manual fuel reporting in Indonesian logistics companies, such as PT Balaraja Distribusindoraya, often leads to inefficiency, fraud, and lack of anomaly supervision. This research aims to develop a web-based system that integrates machine learning and computer vision to monitor fuel consumption and detect anomalies in logistics fleets. The proposed system employs Isolation Forest for unsupervised anomaly detection based on fuel volume, travel distance, and fuel ratio, combined with a deep learning–based CBIR module using MobileNetV2 to validate fuel station images, and OCR to extract numerical data from receipts. Following the CRISP-DM methodology, the model was trained and deployed through a Flask-based API and evaluated using black-box and white-box testing. Experimental results show that Isolation Forest achieves the highest anomaly detection performance (F1-Score = 0.81, ROC-AUC = 0.99), CBIR validates official fuel stations with ≥95% similarity, and OCR reaches 97% accuracy in receipt recognition. The novelty of this study lies in its hybrid integration of anomaly detection and visual verification within a single scalable platform. This research contributes to Informatics by providing a framework for hybrid anomaly detection systems that enhance digitalization, transparency, and operational efficiency in the logistics sector.

*Keywords :* *Anomaly Detection, CBIR, Isolation Forest, Machine Learning, OCR.*

## 1. INTRODUCTION

The use of Machine Learning (ML) technology has been expanding across various industrial sectors, including logistics and transportation, due to its ability to analyze large-scale data and provide accurate predictions and anomaly detection [1]. According to Su et al. (2023), the application of ML technology in logistics management offers deep insights into vehicle fuel consumption, enabling route optimization and identifying irregular or anomalous consumption patterns that may indicate waste or fraud [2]. In the context of fleet management, ML can be utilized to analyze fuel consumption patterns, travel distances, and driver behavior with the aim of improving operational efficiency and reducing fuel costs. However, in Indonesia, most logistics companies still rely on manual reporting systems, which are prone to fraud and inefficiency, thus creating an urgent need for digital-based anomaly detection and monitoring solutions.

The concept of Intelligent Fleet Fuel Management (IFFM) has emerged as a solution for integrating telematics data, vehicle sensors, and GPS to monitor fuel consumption in real time, one example is research conducted by Barbado et al. in 2022 [3]. In practice, this technology not only predicts fuel requirements but also detects fuel usage anomalies through unsupervised learning approaches capable of identifying abnormal data patterns without requiring labeled datasets [4]. One algorithm proven to be effective is Isolation Forest, as it can detect outlier data with low complexity in

large datasets [5]. Recent work on fleet vehicle data anomaly detection, such as the integration of statistical thresholding and Model-in-the-Loop calibration for OBD systems (Kumar et al., 2024), highlights the relevance of real-time automated anomaly monitoring in fleet operations [6]. Nevertheless, most existing IFFM solutions focus only on anomaly detection from numerical or telematics data, without integrating visual verification that could prevent manipulation of transaction evidence.

The use of unsupervised learning approaches for anomaly detection in fuel consumption data has become a prominent focus of research in recent years [7]. Atemkeng et al. (2023) developed a Label Assisted Autoencoder to detect anomalies in fuel consumption for power plants, demonstrating high accuracy even under limited label conditions a method that can be adapted for logistics vehicles [8]. In the context of fleet management, Giannoulidis et al. (2023) proposed a context-aware predictive maintenance solution based on unsupervised methods such as TranAD and clustering techniques, which proved effective in detecting anomalies and supporting deployment in real operational environments [9]. Meanwhile, Muhammad et al. (2024) integrated a Gaussian Mixture Model (GMM) with an Ensemble Isolation Forest in a vehicle telemetry system to identify driving behaviors that affect emissions and fuel consumption, thereby supporting the data understanding and modeling stages within the CRISP-DM framework [10]. Although these studies reported promising accuracy, none of them addressed the issue of validating the authenticity of fuel station photos or transaction receipts, which is crucial in preventing fraud in logistics operations.

In the logistics domain, the ensemble hybrid unsupervised anomaly detection approach developed by Phiboonbanakit et al. (2019) demonstrated improved accuracy in detecting issues within fleet management systems, which is methodologically relevant for optimizing fuel consumption [11]. Finally, Jesmeen et al. (2021) proposed a clustering- and Principal Component Analysis (PCA)-based model for anomaly detection in energy consumption, which can be directly adapted to address fuel consumption in logistics vehicles [12]. However, these approaches primarily concentrate on anomaly detection performance, with little emphasis on integration with computer vision modules such as CBIR or OCR. As a result, their applicability to contexts with high risks of visual manipulation, such as Indonesian logistics operations, remains limited.

In Indonesia, fuel monitoring in logistics companies is still predominantly performed manually, such as recording via WhatsApp and Google Spreadsheet, which is prone to errors and fraud, as exemplified by PT Balaraja Distribusindoraya. Cases of transaction receipt manipulation and refueling at unauthorized fuel stations often create opportunities for waste and operational losses. To address these challenges, a system is needed that not only detects anomalies in fuel consumption but also verifies data authenticity using deep learning–based Content-Based Image Retrieval (CBIR) with MobileNetV2 to check the validity of fuel station images, as well as Optical Character Recognition (OCR) to extract numerical data from fuel receipts. This combination ensures that anomaly detection results are supported by visual evidence validation, reducing the possibility of fraudulent reporting.

Several machine learning algorithms, such as Random Forest, Gradient Boosting, Support Vector Machine (SVM), and Neural Networks, have been employed for fuel consumption analysis [13], [14], [15], [16]. However, research focusing on fuel consumption anomaly detection with the integration of web-based CBIR and OCR remains relatively scarce, particularly in the context of logistics companies in Indonesia. Existing works mostly treat anomaly detection and visual verification as separate problems, while this study integrates them into a unified system, representing a clear research gap.

Therefore, this study aims to develop an optimized fuel monitoring system for logistics companies using the Isolation Forest algorithm for anomaly detection, integrated with a CBIR module based on MobileNetV2 for visual verification and OCR for transaction data extraction, and implemented within an integrated web application. The novelty of this research lies in its hybrid framework that combines

unsupervised anomaly detection with computer vision–based validation, which to the best of our knowledge has not been previously implemented in Indonesian logistics. The contribution to Informatics/Computer Science is the provision of a scalable framework for hybrid anomaly detection that enhances transparency, digitalization, and operational efficiency in fleet management. To ensure optimal performance, the model is compared with DBSCAN and One-Class SVM, subjected to hyperparameter tuning such as adjusting the number of trees (n_estimators) and the maximum tree depth in Isolation Forest and evaluated using Precision, Recall, F1-Score, ROC-AUC, and anomaly score analysis. The expected outcome is an accurate anomaly detection model and a web application capable of helping companies reduce fuel waste, detect potential fraud, and support the digitalization of fleet fuel management in a smarter and more sustainable manner.

## 2. METHOD

This study adopts a quantitative experimental approach by applying a machine learning algorithm based on unsupervised learning to detect anomalies in fuel consumption for logistics vehicles. The model was developed following the CRISP-DM (Cross-Industry Standard Process for Data Mining) methodology, which consists of six stages: business understanding, data understanding, data preparation, modeling, evaluation, and deployment [10].

### 2.1. Research Data

The success of anomaly detection and visual verification in this study largely depends on the quality and relevance of the dataset used. Therefore, the research employed real operational data obtained directly from the company to ensure that the developed system reflects actual field conditions. The dataset consists of transaction records, vehicle usage information, and supporting images that represent the core variables influencing fuel consumption monitoring. Details of the collected data are as follows:

1) Data Source: The data were obtained from the fuel transaction records of 14 operational vehicles at PT Balaraja Distribusindoraya (Wings Group) in Tangerang Regency, covering the period from October 2023 to September 2024.
2) Data Features: These include fuel volume (liters), travel distance (km), fuel consumption ratio (km/liter), vehicle type, and photographs of the fuel stations where refueling took place.
3) Additional Data: A set of four official fuel stations (SPBU 3442123, 3442117, 3442120, 3442107) was used as the validation dataset for the CBIR system
4) Dataset Size: In total, more than 2,000 transaction records were collected, with 70% allocated for training and 30% for testing, ensuring the evaluation process was statistically valid.

### 2.2. Data Processing Techniques

Before the data could be utilized for modeling and system development, several preprocessing steps were necessary to ensure its accuracy, consistency, and suitability for machine learning analysis. Since the raw dataset contained incomplete entries, redundant records, and variations in image quality, it was important to apply a structured data processing workflow. These steps not only improved the reliability of the dataset but also enhanced the model's ability to detect anomalies and validate refueling activities effectively.

The data processing techniques implemented in this study include the following:

1) Data Cleaning: Removing duplicate entries, filling in missing values, and normalizing numerical features.
2) Feature Engineering: Calculating the fuel consumption ratio per vehicle, encoding vehicle types, and adding correlation features for pattern analysis.

3) Image Preprocessing: Processing fuel station photographs using MobileNetV2 for feature vector extraction, followed by cosine similarity measurement to validate refueling locations.
4) OCR Processing: Extracting transaction details (fuel volume, price per liter, total cost) from receipts using Optical Character Recognition to cross-validate with recorded data.

### 2.3. Machine Learning Model

After the dataset was cleaned and preprocessed, the next stage focused on developing machine learning models for anomaly detection in fuel consumption. The choice of algorithm plays a critical role in ensuring the system's ability to distinguish between normal and anomalous patterns effectively. In this study, the Isolation Forest algorithm was selected as the primary model due to its efficiency and robustness in handling large-scale, unlabeled data. To provide a comprehensive comparison and validate the reliability of the approach, additional models such as DBSCAN and One-Class SVM were also implemented as benchmarks. The modeling process involved hyperparameter tuning, evaluation using multiple performance metrics, and integration with other system modules to achieve optimal results.

1) Primary Algorithm: Isolation Forest was selected as the main anomaly detection model due to its efficiency in isolating outlier data with low computational complexity. The algorithm was configured with n_estimators = 100, contamination = 0.1, and max_samples = 256. The anomaly score was calculated using Equation (1):

$$s(x,n) = 2 - \frac{E(h(x))}{c(n)} \qquad (1)$$

where $E(h(x))$ represents the expected path length of instance $x$ and $c(n)$ is the average path length of unsuccessful searches in a binary tree of size $n$.

2) Comparison Models: DBSCAN and One-Class SVM were used for performance benchmarking.
3) Hyperparameter Tuning: Optimization of hyperparameters was performed, including the number of trees (n_estimators) and the maximum tree depth for the Isolation Forest model.

### 2.4. Model Evaluation

Once the machine learning models were developed, it was essential to evaluate their performance to ensure reliability and practical applicability. Model evaluation not only validates the accuracy of anomaly detection but also determines how well each algorithm can generalize to unseen data. In this study, multiple evaluation metrics were employed, including Precision, Recall, F1-Score, and ROC-AUC, as these provide a balanced perspective on both detection capability and error minimization. Furthermore, manual validation using domain knowledge was applied to strengthen the ground truth, ensuring that the identified anomalies were contextually relevant to actual fuel consumption behavior. In addition to numerical evaluation, the CBIR module was assessed using cosine similarity thresholds (0.8 as acceptance level) and Top-k retrieval accuracy. OCR performance was measured by calculating character-level accuracy and word-level recognition rates on extracted receipt data. This comprehensive evaluation process allowed for an objective comparison between Isolation Forest, DBSCAN, and One-Class SVM, and guided the selection of the most effective model for integration into the system.

### 2.5. System Integration

After the best-performing model was identified through the evaluation stage, the next step was to integrate the anomaly detection and validation modules into a functional system. This integration was carried out to ensure that the developed machine learning models could be accessed and utilized effectively by end users through a web-based application. The process involved embedding the trained model into a Flask-based API, connecting it with the application's backend, and designing user-friendly interfaces for administrators and drivers. In addition to anomaly detection, supporting modules such as

CBIR for gas station image validation and OCR for receipt data extraction were also combined into the system. The system architecture was designed in a modular way, allowing for future scalability such as integration with IoT sensors or GPS-based telematics data. A research flow diagram (Figure 1) was also included to illustrate the workflow from data collection, preprocessing, modeling, evaluation, and deployment.
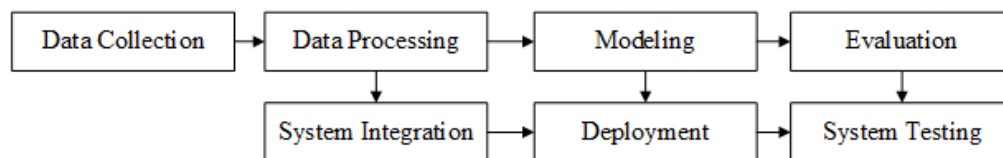


Figure 1. Hybrid Anomaly Detection Framework

## 2.6. System Testing

After the anomaly detection model and supporting modules were successfully integrated into the web-based application, the next step was to evaluate the system as a whole to ensure that it operated according to the intended design. System testing was conducted using two complementary approaches, namely black-box and white-box testing. Black-box testing focused on verifying the functional aspects of the application from the user's perspective, such as login, data input, anomaly prediction, and validation of receipts and fuel station photographs. Meanwhile, white-box testing was carried out to examine the internal logic of the algorithms, the accuracy of anomaly predictions, and the stability of API communication between the backend and frontend. This dual testing was supplemented by stress testing with simulated high-volume data input (up to 500 transactions per batch) to evaluate system scalability and performance stability under load conditions.

## 3. RESULT

This section presents the experimental results organized according to the CRISP-DM stages described in the methodology, namely business understanding, data understanding, data preparation, modeling, evaluation, and deployment. This structure ensures consistency between the proposed method and the reported results while clearly demonstrating the effectiveness of the developed hybrid anomaly detection and visual validation system.

## 3.1. Business Understanding Result

Based on field observations at PT Balaraja Distribusindoraya, manual fuel monitoring practices using WhatsApp messages and Google Spreadsheets were identified as the primary source of inefficiency and fraud risk. Common issues included delayed reporting, inconsistent records, manipulation of fuel receipts, and refueling at unauthorized gas stations. These findings confirm the necessity of an automated and intelligent fuel monitoring system capable of detecting consumption anomalies and validating transaction authenticity in real operational conditions.

## 3.2. Data Understanding Result

The collected dataset consisted of 4,213 fuel transaction records with 26 attributes, obtained from 14 operational vehicles over a one-year period. Initial exploratory analysis revealed variations in fuel consumption ratios across vehicle types and operational routes. Statistical analysis was conducted to examine minimum, maximum, median, and average values, as well as missing entries. Feature correlation analysis showed strong relationships between fuel consumption ratio and vehicle type, indicating their importance for anomaly detection. Figure 1 illustrates the correlation heatmap used to support feature selection.
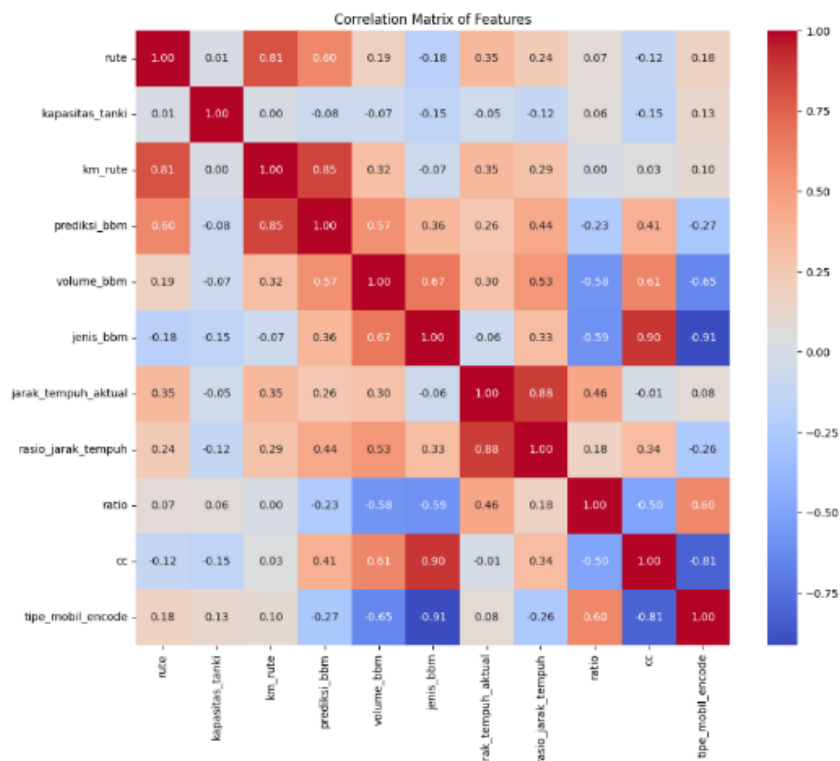
Figure 1. Feature Correlation Visualization

### 3.3. Data Preparation Result

During data preparation, irrelevant attributes were removed, missing values were filled using median or average imputation, and numerical features were normalized. Vehicle type data were encoded and ranked based on fuel consumption ratios. The cleaned dataset enabled stable modeling and reduced noise caused by incomplete or redundant records. The resulting structured dataset served as input for both anomaly detection models and visual validation modules.

### 3.4. Modeling Result

Three unsupervised learning algorithms were applied to detect anomalies in fuel consumption patterns: Isolation Forest, DBSCAN, and One-Class SVM. Among them, Isolation Forest demonstrated the most consistent separation between normal and anomalous data points, as visualized in Figure 2. Anomaly thresholds were refined through manual validation using domain knowledge to ensure contextual relevance.
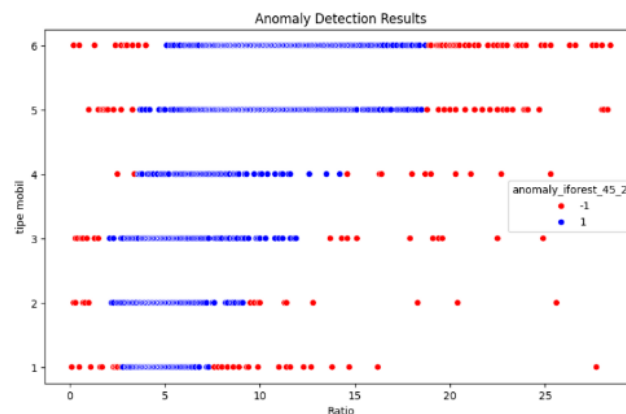


Figure 2. Isolation Forest Model Results

DBSCAN showed a looser clustering behavior, retaining more extreme values as normal (Figure 3), while One-Class SVM classified a larger portion of the dataset as anomalous, potentially increasing false positives (Figure 4).
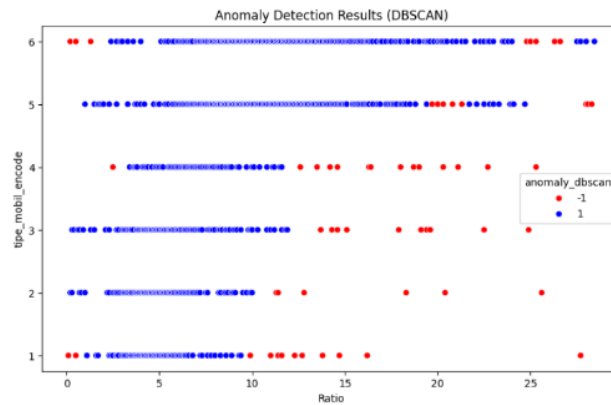


Figure 3. DBSCAN Model Results



Figure 4. One-Class SVM Model Results

### 3.5.  Evaluation Result

The model evaluation was conducted through manual identification of the anomaly scores generated by the model. The manually identified results were then used as the ground truth, serving as domain knowledge for determining whether a data point was anomalous. Subsequently, the model predictions were compared with the manual validation results derived from domain knowledge, and the evaluation metrics Precision, Recall, F1-Score, and ROC-AUC were calculated.

The evaluation results indicate that Isolation Forest delivered the best performance in detecting anomalies in fuel transaction data. Table 1 shows that Isolation Forest outperformed DBSCAN and One-Class SVM, achieving the highest F1-Score (0.816) and ROC-AUC (0.995), indicating robust anomaly detection performance.

Table 1.  Model Score Comparison

| Model | Precision | Recall | F1-Score | ROC-AUC |
|---|---|---|---|---|
| Isolation Forest | 0.807947 | 0.824324 | 0.816054 | 0.995076 |
| DBSCAN | 0.741667 | 0.601351 | 0.664179 | NaN |
| One-Class SVM | 0.373134 | 0.675676 | 0.480769 | 0.860763 |

Figure 5 compares the overall performance of the three models, confirming the superiority of Isolation Forest.
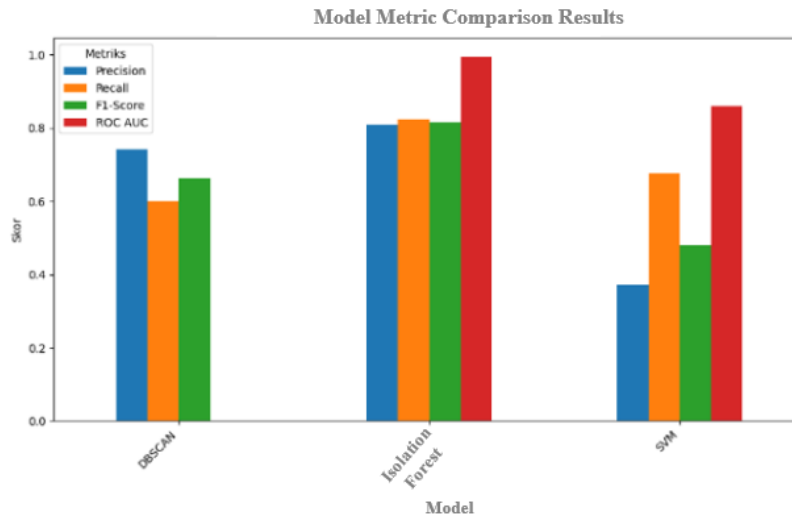


Figure 5. Model Comparison (DBSCAN, Isolation Forest, and One-Class SVM)

As shown in Table 1 and Figure 5, Isolation Forest not only outperformed the other models in terms of F1-Score but also demonstrated stability in the ROC-AUC evaluation.

Furthermore, the hyperparameter tuning process for Isolation Forest revealed that setting n_estimators to 100 and contamination to 0.042 yielded optimal performance (Figure 6). The ROC curve and precision-recall plot illustrate a well-balanced trade-off between anomaly detection accuracy and minimizing false positives.

```python
from sklearn.ensemble import IsolationForest
from sklearn.model_selection import GridSearchCV

# Model Isolation Forest
model = IsolationForest(random_state=42)

# Rentang hyperparameter untuk dicoba
param_grid = {
    'n_estimators': [100],
    'max_samples': ['auto'],
    'contamination': [0.04,0.05, 0.1, 0.2, 0.3],
    'max_features': [1, 2]
}

# Grid Search
grid_search = GridSearchCV(estimator=model, param_grid=param_grid, cv=5, scoring='accuracy')
grid_search.fit(X_train, y_train)

# Hasil terbaik
print("Best Hyperparameters:", grid_search.best_params_)
print("Best Score:", grid_search.best_score_)
```

```
Best Hyperparameters: {'contamination': 0.04, 'max_features': 2, 'max_samples': 'auto', 'n_estimators': 100}
Best Score: 0.8998430870953988
```

Figure 6. Hyperparameter Tuning

### 3.6.  Visual Validation Result

To prevent fraudulent fuel refills at unauthorized stations, the system is equipped with a MobileNetV2-based CBIR (Content-Based Image Retrieval) module. Test results indicate that the system can distinguish fuel station images with a similarity level of $\geq 95\%$ using cosine similarity. Accordingly, every fuel receipt photo and fuel station image uploaded by drivers is automatically validated to ensure compliance with the official stations (SPBU 3442123, 3442117, 3442120, 3442107). Meanwhile, the OCR (Optical Character Recognition) module used to extract numerical values from fuel receipts achieved an accuracy rate of 97% after image preprocessing steps, including binarization,

brightness adjustment, and sharpening. This combination further enhances the visual validation system for fuel refilling reports. An illustration of the activity diagram is shown below (Figure 7).
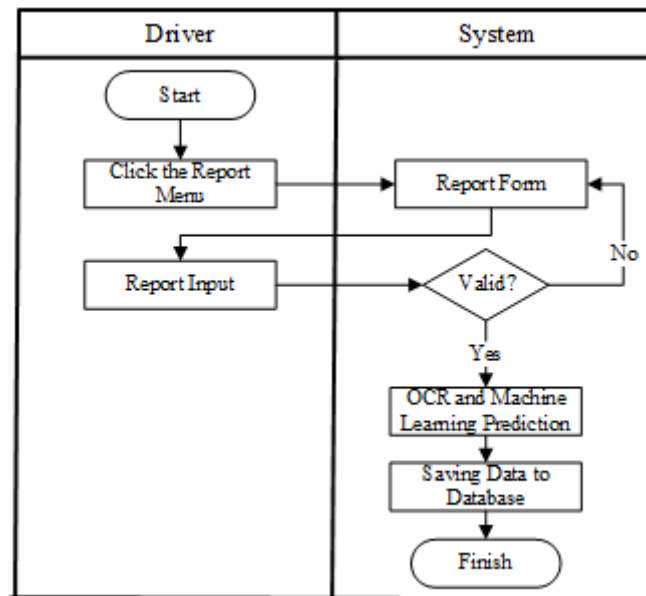


Figure 7. Activity Diagram for Fuel Report Input

### 3.7. Deployment and System Testing

The Isolation Forest model was integrated into the backend system via a Flask-based API. The frontend was developed using Laravel and Bootstrap, supported by a MySQL database. The application provides two main interfaces (Figure 8):

1) Admin Dashboard – Used to monitor fuel transaction reports, view anomaly detection results, and manage vehicle and driver data.
2) Driver Interface – Used to upload fuel reports, fuel station photos, and fuel receipt images.
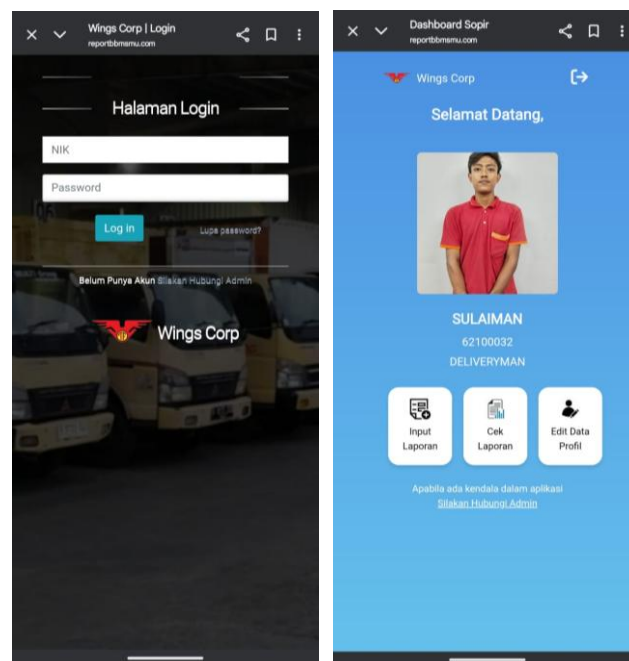


Figure 8. User Interface (Deliveryman)

Black-box testing of the application showed that all core functions operated according to specifications, while white-box testing verified that the anomaly detection logic and image validation processes worked as intended (Table 2).

Table 2. Summary of System Testing Results

| Component | Result |
| --- | --- |
| Login and Authentication | Successful |
| Fuel Report Input | Successful |
| Anomaly Prediction | Accurate (F1=0.81) |
| OCR Validation | 97% Acurate |
| CBIR Validation | ≥95% Match |
| API Response | Stable |

In addition to the summarized results in Table 2, several representative scenarios were conducted to validate the system in practical use cases. For instance, when a driver uploaded a photo of an unauthorized fuel station, the CBIR module successfully rejected it with a similarity score below 60%, marking the transaction as suspicious. In another scenario, an uploaded fuel receipt with blurred text was still correctly processed by the OCR module, achieving 96% accuracy after preprocessing steps such as brightness adjustment and sharpening. Furthermore, the Isolation Forest model effectively detected an anomalous report where a vehicle recorded 5 liters of fuel but was logged to travel 350 km, classifying it as abnormal with a high anomaly score. These scenarios confirm that the system functions reliably under both typical and irregular conditions, thereby strengthening its applicability in real-world operations. The system implementation was also tested directly by internal users at PT Balaraja Distribusindoraya, receiving positive feedback for its ease of use and detection speed.

These results indicate that the proposed system not only achieves high accuracy in anomaly detection but also delivers practical benefits for logistics operations. By automating data validation through CBIR and OCR, as well as integrating anomaly detection within a single web-based platform, the system reduces the potential for human error and fraud, improves operational transparency, and accelerates decision-making processes for fleet management.

## 4. DISCUSSIONS

Based on the testing results, the Isolation Forest algorithm demonstrated the best performance in detecting anomalies in fuel consumption compared to DBSCAN and One-Class SVM. This finding is consistent with the studies by Xiang et al. (2023) and Utkin et al. (2023), which confirmed the effectiveness of Isolation Forest in isolating outliers through a low-complexity isolation tree mechanism, as well as its ability to maintain performance stability across large data variations [17], [18]. In this study, the optimal configuration was achieved with n_estimators = 100 and contamination = 0.042, resulting in an F1-Score of 0.81 and a ROC-AUC score close to 0.99.

When compared to DBSCAN, although this density-based method can effectively cluster data, it tends to retain more extreme data points as "normal," which leads to a lower recall rate (0.60). One-Class SVM, on the other hand, is more sensitive to data variations and produces a higher false positive rate, despite achieving a relatively competitive recall (0.67). These results are in line with the findings

of Hurst et al. (2020) and Airlangga (2024), which highlight the limitations of DBSCAN and One-Class SVM in detecting anomalies within datasets that have heterogeneous distributions [19], [20] .

Table 3.  System Model & Module Comparison Summary

| Component/Model | Key Advantages | Limitations / Challenges | Results in This Study |
|---|---|---|---|
| Isolation Forest | ▪ Effective anomaly detection in large, unlabeled datasets (unsupervised). ▪ Low complexity and stable across data variations. | ▪ Requires parameter tuning for optimal performance | ▪ F1-Score: 0.81. ▪ ROC-AUC: 0.99. ▪ Optimal configuration: n_estimators = 100, contamination = 0.042 |
| DBSCAN | ▪ Detects patterns based on density without labels. ▪ Capable of finding clusters with arbitrary shapes. | ▪ Sensitive to eps and minPts parameters. ▪ Tends to classify extreme data points as normal | ▪ Precision: 0.74 ▪ Recall: 0.60 |
| One-Class SVM | ▪ Separates normal and anomalous data non-linearly using kernels | ▪ Prone to high false positive rates. ▪ More computationally intensive | ▪ Precision: 0.37. ▪ Recall: 0.67 |
| CBIR (MobileNetV2) | ▪ Transfer learning enables high accuracy on resource-constrained devices. ▪ Effective in recognizing images with complex visual features. | Requires a sufficiently representative dataset of official photographs | Similarity rate ≥ 95% for official gas station photos |
| OCR (PP-OCR + Preprocessing) | ▪ High accuracy in reading text from images. ▪ Preprocessing improves robustness for blurry or dark images. | Sensitive to heavy noise or extremely skewed photos | Accuracy: 97% in reading fuel receipt data. |

For CBIR validation, using MobileNetV2 as the feature extractor yielded a recognition accuracy of ≥ 95% for official gas station images, based on cosine similarity. This approach adopted a transfer learning strategy, enabling the system to run efficiently while maintaining accuracy on resource-constrained devices, as demonstrated by Ahmad Saeed Mohammad et al. (2024). This is particularly

relevant given that one common fraud scheme in the field involves fueling at unofficial stations or manipulating photographic transaction evidence [21].

Meanwhile, the OCR module used for reading fuel receipt data achieved a 97% accuracy rate after image preprocessing (binarization, brightness adjustment, and sharpening). This result supports the findings of Chenxia Li et al. (2022) regarding the effectiveness of PP-OCRv3 and optimized text recognition techniques in low-quality image conditions. The integration of OCR significantly streamlines the data input process and minimizes the risk of manual entry errors [22].

From a practical standpoint, the combination of Isolation Forest-based anomaly detection, CBIR visual validation, and OCR provides a fuel monitoring system that not only detects irregular consumption but also verifies the authenticity of transaction evidence. This makes the system a relevant solution for logistics companies in Indonesia, which often face challenges related to manual reporting and low transparency in supervision.

From a business perspective, the implementation of this system delivers a direct impact on the operational efficiency and transparency of logistics companies [23], [24]. With a high anomaly detection accuracy (F1-Score: 0.81) and automated transaction evidence validation [23], companies can identify potential fuel fraud or wastage more quickly, enabling them to mitigate losses before they become significant. Furthermore, the time required to verify fuel refueling reports is drastically reduced, as validation is performed automatically through CBIR and OCR technologies [25]. This reduces the manual workload of administrative staff and allows management teams to make data-driven decisions more rapidly [26].

In terms of transparency, the system provides a clear digital trail for every fuel transaction, minimizing opportunities for data manipulation and enhancing accountability for both drivers and partner gas stations [27]. Operational cost efficiency is also improved through the reduction of irregular fuel consumption and better vehicle maintenance planning based on validated consumption data. Consequently, this research not only contributes technically to the application of machine learning and computer vision but also delivers a strategic impact on the efficiency and integrity of company operations (in Table 3).

## 5. CONCLUSION

This study developed a web-based system to monitor and detect anomalies in fuel consumption for the operational vehicles of PT Balaraja Distribusindoraya. The system integrates the Isolation Forest algorithm for anomaly detection, a CBIR module based on MobileNetV2 for gas station photo verification, and OCR for reading data from fuel receipts. The development process followed the CRISP-DM methodology and was implemented through a Flask-based API connected to a web application. The evaluation results show that:

1) Isolation Forest outperforms DBSCAN and One-Class SVM, achieving an F1-Score of 0.81 and an ROC-AUC of 0.99.
2) The CBIR module successfully recognizes official gas station photos with a similarity score of $\geq$ 95%, while OCR achieves 97% accuracy in reading receipt data.
3) The system performs reliably in both black-box and white-box testing, and received positive feedback from the company's internal users.

The contribution of this research lies in providing an integrated solution that enhances the accuracy of fuel anomaly detection, minimizes the potential for fraud, and supports the digitalization of fleet monitoring processes. Future developments could focus on adding real-time visualization features, integrating additional variables (such as route type and driver behavior), applying ensemble or deep learning methods, and implementing automatic notifications to accelerate responses to anomalies.

## ACKNOWLEDGEMENT

## REFERENCES

[1] K. Tsolaki, T. Vafeiadis, A. Nizamis, D. Ioannidis, and D. Tzovaras, "Utilizing Machine Learning on Freight Transportation and Logistics Applications: A review," *ICT Express*, pp. 284–295, June 2023, doi: 10.1016/j.icte.2022.02.001.

[2] M. Su, Z. Su, S. Cao, K.-S. Park, and S.-H. Bae, "Fuel Consumption Prediction and Optimization Model for Pure Car/Truck Transport Ships," *J. Mar. Sci. Eng.*, vol. 11, no. 6, June 2023, doi: 10.3390/jmse11061231.

[3] A. Barbado and Ó. Corcho, "Interpretable machine learning models for predicting and explaining vehicle fuel consumption anomalies," *Eng. Appl. Artif. Intell.*, vol. 115, p. 105222, Oct. 2022, doi: 10.1016/j.engappai.2022.105222.

[4] B. B. Turan, E. Genç, İ. N. Akçığ, N. Göztepe, M. E. Mumcuoğlu, and M. Ünel, "Detecting high fuel consumption in HDVs with ensemble of anomaly detection models," Aug. 2024, Accessed: July 29, 2025. [Online]. Available: http://dx.doi.org/10.1109/INDIN58382.2024.10774415

[5] H. Xiang *et al.*, "OptIForest: Optimal Isolation Forest for Anomaly Detection," June 23, 2023, *arXiv*: arXiv:2306.12703. doi: 10.48550/arXiv.2306.12703.

[6] A. Kumar, K. Hegde, K. Challa, and Y. H, "Anomaly Detection in Fleet Vehicle Data and Statistical Approach to Develop Engine System OBD Thresholds," *SAE Publ.*, p. 10, Dec. 2024, doi: https://doi.org/10.4271/2024-28-0195.

[7] O. F. M. ElMahdy, M. E. Hassan, and S. M. Metwalli, "Machine learning anomaly detection of lost and unaccounted for gas in natural gas networks," *J. Eng. Appl. Sci.*, vol. 72, no. 1, p. 123, Aug. 2025, doi: 10.1186/s44147-025-00677-x.

[8] M. Atemkeng *et al.*, "Label Assisted Autoencoder for Anomaly Detection in Power Generation Plants," Feb. 06, 2023, *arXiv*: arXiv:2302.02896. doi: 10.48550/arXiv.2302.02896.

[9] A. Giannoulidis and A. Gounaris, "A context-aware unsupervised predictive maintenance solution for fleet management," *J. Intell. Inf. Syst.*, vol. 60, no. 2, pp. 521–547, Apr. 2023, doi: 10.1007/s10844-022-00744-2.

[10] A. S. Muhammad, C. Wang, and L. Chen, "A Telemetric Framework for Assessing Vehicle Emissions Based on Driving Behavior Using Unsupervised Learning," *Vehicles*, vol. 6, no. 4, pp. 2170–2194, Dec. 2024, doi: 10.3390/vehicles6040106.

[11] T. Phiboonbanakit, V.-N. Huynh, T. Horanont, and T. Supnithi, "Unsupervised hybrid anomaly detection model for logistics fleet management systems," *IET Intell. Transp. Syst.*, vol. 13, no. 11, Nov. 2019, Accessed: Aug. 16, 2025. [Online]. Available: https://trid.trb.org/View/1697032

[12] J. M. Z. H, J. Hossen, and A. B. A. Aziz, "Unsupervised Anomaly Detection for Energy Consumption in Time Series using Clustering Approach," *Emerg. Sci. J.*, vol. 5, no. 6, pp. 840–854, Dec. 2021, doi: 10.28991/esj-2021-01314.

[13] N. Norouzkhani *et al.*, "Developing and evaluating a gamified self-management application for inflammatory bowel disease using the ADDIE model and Sukr framework," *BMC Med. Inform. Decis. Mak.*, vol. 25, no. 1, p. 11, Jan. 2025, doi: 10.1186/s12911-024-02842-3.

[14] G. M. Nguegnang, M. Atemkeng, T. Ansah-Narh, R. Rockefeller, G. M. Nguegnang, and M. A. Garuti, "Predicting Fuel Consumption in Power Generation Plants using Machine Learning and Neural Networks," Feb. 11, 2022, *arXiv*: arXiv:2202.05591. doi: 10.48550/arXiv.2202.05591.

[15] D. Zhao *et al.*, "A Review of the Data-Driven Prediction Method of Vehicle Fuel Consumption," *Energies*, vol. 16, no. 14, Art. no. 14, Jan. 2023, doi: 10.3390/en16145258.

[16] A. Fan, Y. Wang, L. Yang, X. Tu, J. Yang, and Y. Shu, "Comprehensive evaluation of machine learning models for predicting ship energy consumption based on onboard sensor data," *Ocean Coast. Manag.*, vol. 248, p. 106946, Feb. 2024, doi: 10.1016/j.ocecoaman.2023.106946.

[17]    H. Xiang *et al.*, "OptIForest: Optimal Isolation Forest for Anomaly Detection," in *Proceedings of the Thirty-Second International Joint Conference on Artificial Intelligence*, Macau, SAR China: International Joint Conferences on Artificial Intelligence Organization, Aug. 2023, pp. 2379–2387. doi: 10.24963/ijcai.2023/264.

[18]    L. Utkin, A. Ageev, A. Konstantinov, and V. Muliukha, "Improved Anomaly Detection by Using the Attention-Based Isolation Forest," *Algorithms*, vol. 16, no. 1, p. 19, Jan. 2023, doi: 10.3390/a16010019.

[19]    W. Hurst, C. A. C. Montañez, and N. Shone, "Time-Pattern Profiling from Smart Meter Data to Detect Outliers in Energy Consumption," *IoT*, vol. 1, no. 1, pp. 92–108, Sept. 2020, doi: 10.3390/iot1010006.

[20]    G. Airlangga, "Advanced Machine Learning Techniques For Seismic Anomaly Detection In Indonesia: A Comparative Study Of LOF, Isolation Forest, And One-Class SVM," *J. Lebesgue J. Ilm. Pendidik. Mat. Mat. Dan Stat.*, vol. 5, no. 1, pp. 49–61, Apr. 2024, doi: 10.46306/lb.v5i1.490.

[21]    A. S. Mohammad, T. G. Jarullah, M. T. S. Al-Kaltakchi, J. Alshehabi Al-Ani, and S. Dey, "IoT-MFaceNet: Internet-of-Things-Based Face Recognition Using MobileNetV2 and FaceNet Deep-Learning Implementations on a Raspberry Pi-400," *J. Low Power Electron. Appl.*, vol. 14, no. 3, p. 46, Sept. 2024, doi: 10.3390/jlpea14030046.

[22]    C. Li *et al.*, "PP-OCRv3: More Attempts for the Improvement of Ultra Lightweight OCR System," June 14, 2022, *arXiv*: arXiv:2206.03001. doi: 10.48550/arXiv.2206.03001.

[23]    F. Moradi, M. Tarif, and M. Homaei, "Semi-Supervised Supply Chain Fraud Detection with Unsupervised Pre-Filtering," Aug. 07, 2025, *arXiv*: arXiv:2508.06574. doi: 10.48550/arXiv.2508.06574.

[24]    H. Wang, L. S. Sua, and B. Alidaee, "Enhancing supply chain security with automated machine learning," July 22, 2025, *arXiv*: arXiv:2406.13166. doi: 10.48550/arXiv.2406.13166.

[25]    K. Sivakoti, "Vehicle Detection and Classification for Toll collection using YOLOv11 and Ensemble OCR," Dec. 13, 2024, *arXiv*: arXiv:2412.12191. doi: 10.48550/arXiv.2412.12191.

[26]    K. Thammarak, Y. Sirisathitkul, P. Kongkla, and S. Intakosum, "Automated Data Digitization System for Vehicle Registration Certificates Using Google Cloud Vision API," *Civ. Eng. J.*, vol. 8, no. 7, pp. 1447–1458, July 2022, doi: 10.28991/CEJ-2022-08-07-09.

[27]    D. Fleischhacker, W. Goederle, and R. Kern, "Improving OCR Quality in 19th Century Historical Documents Using a Combined Machine Learning Based Approach," Jan. 15, 2024, *arXiv*: arXiv:2401.07787. doi: 10.48550/arXiv.2401.07787.