

# Final Project CSE527: The Interpretation of Vanity License Plates

## Abstract

Vanity license plates enable drivers to express individuality but pose safety and administrative challenges. Currently, DMV offices rely on manual reviews to evaluate vanity plate applications, screening for offensive language and regional slang. This process is labor-intensive, costly, and often struggles to interpret diverse linguistic and cultural nuances, resulting in inappropriate approvals or unjust rejections.

This study proposes an automated framework to enhance vanity plate management efficiency and compliance. The system incorporates:

- Personalized Plate Detection: A Lexicon-based method to distinguish vanity plates from standard ones.
- Semantic Interpretation: Large language models to decode the meaning of vanity plates, including slang and cultural references.
- Content Compliance Assessment: Automated identification of offensive or illegal content to support DMV evaluations.

The proposed solution aims to streamline application reviews, mitigate inappropriate approvals, and provide a scalable tool for social psychologists exploring the behavioral implications of vanity plates.

## 1. Introduction

With the rapid proliferation of personalized vehicle identification systems, vanity license plates have become a popular means for drivers to express individuality[2]. However, this trend presents significant challenges for effective management and oversight. Current manual and rule-based evaluation processes employed by Departments of Motor Vehicles (DMVs) struggle to keep pace with the growing volume of vanity plate applications, which often involve creative abbreviations, diverse linguistic expressions, and cultural nuances. These inefficiencies can lead to inappropriate approvals, unfair rejections, or processing delays, undermining administrative efficiency and public trust[18].

Beyond administrative challenges, vanity plates also pose safety risks. Studies have highlighted that visual distractions, including those caused by vanity plates, contribute to up to 83% of traffic incidents, emphasizing the importance of effective screening and regulation[12]. Despite the growing prevalence of personalized plates, automated tools for managing and evaluating these systems remain limited. Existing approaches primarily rely on manual reviews or dictionary-based methods. Manual reviews are labor-intensive and inconsistent, while dictionary-based filters fail to capture evolving slang, implicit meanings, and cultural contexts, resulting in low accuracy and high error rates.

To address these challenges, we propose a novel framework for automating the screening of vanity license plates, focusing on leveraging large language models (LLMs) to detect illegal or offensive content. Unlike traditional dictionary-based methods, LLMs provide the capability to understand nuanced language, cultural

references, and implicit meanings, significantly enhancing the accuracy and efficiency of the evaluation process.

In this study, we use a publicly available dataset of vanity license plates annotated for compliance and legality. The dataset includes diverse examples of valid and invalid plates, providing a robust foundation for developing and testing our model. We fine-tune an LLM on this dataset to detect inappropriate or illegal content and evaluate its performance across various linguistic and cultural contexts.

Comprehensive experiments demonstrate that our system significantly outperforms traditional approaches in both detection accuracy and processing efficiency. The results highlight the potential of LLM-based solutions in transforming the management of vanity plates, reducing the reliance on manual reviews, and improving overall administrative effectiveness.

This work provides a scalable and efficient solution for DMVs, contributing to safer roads and more effective administrative practices. It also sets the stage for future research on applying AI-driven approaches to personalized vehicle identification systems.

## 2. Related Work

In the field of hate speech detection, traditional methods have primarily relied on rule-based systems and machine learning models, such as support vector machines (SVM) and logistic regression, often using handcrafted features like n-grams, word embeddings, and lexical patterns. These methods are effective when specific keywords or phrases related to hate speech are predefined. For instance, Davidson et al. (2017) applied a rule-based approach to detect hate speech on social media platforms, focusing on offensive language and explicit terms like racial slurs or discriminatory expressions. However, these traditional approaches often struggle with the context and subtleties of language, such as sarcasm or implicit hate speech[16].

Recent advancements in Large Language Models (LLMs), such as GPT, BERT, and Llama3, have significantly improved hate speech detection by leveraging large textual corpora to understand context and subtle linguistic cues. Zhang et al. (2020) demonstrated that fine-tuning BERT for hate speech detection allows contextualized models to capture complex relationships in language, outperforming traditional methods[17]. Additionally, Sen et al. (2024) showed that fine-tuning small parameter LLMs with LoRA significantly enhances the model's ability to process hate speech, achieving an accuracy rate of over 80%[11].

However, in the context of personalized license plate evaluation, several challenges arise. Unlike hate speech detection, where context and explicit linguistic cues are often present, license plate content lacks sufficient context and often contains significant textual distortions (e.g., obfuscations, abbreviations). The absence of background information and the disruptive nature of personalized plates introduce unique challenges. Nevertheless, leveraging LLMs remains a promising approach to reduce human effort in screening large volumes of license plate requests.

## 3. Data Collection

### 3.1 Dataset resource

The datasets used in this study were obtained from publicly available sources on GitHub, specifically focused on vanity license plates from California and New York[5,6]. The California dataset comprises 23,463 records collected between 2015 and 2016, including 18,757 rejected plates and 4,673 accepted plates. Each record contains the plate number, reviewer comments, customer-provided explanations, and the approval or rejection status. The New York dataset consists of 133,636 records spanning 2010 to 2014, with 131,990 approved plates, 1,646 rejected plates, and 1,801 flagged as "red\_guide." The dataset provides the plate number and its final classification status.

### 3.2 Dataset Privacy

It is important to note that while the datasets are publicly available, privacy concerns regarding license plate data need to be addressed. The data collected does not include personally identifiable information (PII), as the focus is on the plate number, status, and relevant reviewer comments. However, as license plates can sometimes be linked to personal vehicles, measures should be taken to ensure that any sensitive information is anonymized and that the use of such data complies with privacy regulations. In this study, no direct linkage to individual owners or vehicles was made, and the data was used solely for academic research purposes. Nonetheless, future work should consider additional privacy safeguards to further ensure the confidentiality and ethical use of license plate data.

### 3.3 Preprocess dataset

To ensure data consistency and suitability for training, the following preprocessing steps were applied:

- 1. Removal of Area Code-Related Plates**

Plates referencing area codes were identified and removed exclusively from the California dataset. Such plates were inconsistently classified, with some being approved and others rejected, likely due to variations in the reviewers' knowledge of California's area codes and their potential implications. Since determining the legality of these plates would require a comprehensive understanding of local area codes and their context, we excluded all records containing references to area codes to avoid ambiguity and ensure consistency in the dataset.

- 2. Dataset Integration**

The California and New York datasets were merged into a single unified dataset to provide comprehensive coverage.

- 3. Conflict Resolution**

A conflict occurred when the same plate number appeared in both the California and New York datasets but had different approval/rejection statuses. For instance, a plate might have been approved in one dataset and rejected in another. To ensure consistency and avoid discrepancies in the data, all such conflicting records were removed from the dataset.

- 4. Class Balancing**

Both datasets exhibited significant class imbalance, with far more approved plates than rejected ones. To address this, we balanced the training set by randomly sampling 1,000 approved and 1,000 rejected

plates from each dataset, yielding 4,000 samples in total. The test set was similarly balanced, with 25 approved and 25 rejected plates from each dataset, totaling 100 samples.

## 5. Taxonomy Labeling

Each plate in the test set was manually annotated with taxonomy labels. These labels categorize the plates into specific categories such as legal, offensive, or culturally relevant based on their content and context.

## 6. Random License Plates

We augmented the dataset by generating 10,000 random plates based on the official license plate formats of California (e.g., "1ABC123") and New York (e.g., "ABC1234"), with 5,000 plates generated for each state.

Table 1:

Plates	California	New York
Accepted	4673	131990
Rejected	18757	3447

*This shows the number of accepted and rejected plates in both the California and New York datasets. As illustrated, both datasets exhibit significant class imbalance, with the California dataset having far more rejected plates (18,757) than accepted plates (4,673). Similarly, the New York dataset contains a much larger number of accepted plates (131,990) compared to rejected plates (3,447). This imbalance highlights the need for data balancing techniques in order to ensure fair and*

*accurate model training.*

# 4. Methodology

## 4.1 Taxonomy

To develop a taxonomy for personalized license plates, this study analyzes customer explanations for approved plates and evaluations for rejected plates, identifying common high-frequency themes and primary reasons for rejection. Through frequency analysis and word cloud generation, we aim to uncover the main motivations behind individuals' applications for personalized license plates, as well as the factors that contribute to the rejection of these applications. Based on the analysis of word clouds and referencing Llama Guard's Safety Risk Taxonomy[3], as well as existing regulations on personalized license plate restrictions[1,5], we designed a two-tier classification system for personalized license plates. Fig[1] illustrates the taxonomy of plates we manually constructed. The word clouds and related analysis are provided in Appendix A.



Fig 1: Manually Constructed Vanity Plates Taxonomy

## 4.2 Lexicon-based Method

The Lexicon-based method involves predefining a list of sensitive words and considering various word variants, such as homophones, symbol substitutions, and other modifications, to determine whether a license plate contains any prohibited terms[8]. We used a dictionary of sensitive words along with their relevant variants[7]. This approach allows for rapid, efficient, and transparent initial screening of license plates. However, the method is limited by the coverage of the dictionary, which may fail to identify new or concealed non-compliant content. Furthermore, the lack of contextual understanding in this approach can lead to false positives or missed detections.

## 4.3 LLM Method

The study utilized the **Llama 3.1-8B-Instruct** model to review personalized license plates through simple role design, task description, and response format requirements in the prompt. The table below outlines the prompts used for approval.

## 4.4 LoRA Fine-tuning

To enhance the performance of the LLM on our specific downstream task, we employed **LoRA (Low-Rank Adaptation)** fine-tuning. LoRA fine-tuning optimizes the model's performance by inserting low-rank adaptation layers into the pre-trained model, while also reducing computational and storage costs[9]. We aim to leverage LoRA fine-tuning to enable the model to better identify and assess the content of license plates, particularly detecting potentially obfuscated violations, thereby improving the accuracy of the approval process.

## 4.5 Chain-of-Thought (CoT)

**Chain-of-Thought (CoT)** is a method for solving problems through step-by-step reasoning, which aims to break down complex tasks into simple, logical steps to help the model make accurate decisions[10]. Traditional classification methods typically rely on direct label prediction, whereas the CoT method progressively analyzes the content of the license plate, taking into account its attributes, implied meanings, and related rules to categorize it appropriately.

## 5. Experiments

### 5.1 Identification of Personalized License Plates

Given that both California and New York issue random license plates following specific patterns, we can filter out plates that adhere to these patterns. In California, the license plate format follows the pattern "1ABC123" (one digit, followed by three letters and three digits), while in New York, the format follows the pattern "ABC1234" (three letters followed by four digits). Plates that match these predefined patterns are classified as random license plates and excluded from the personalized plate analysis.

### 5.2 Legality Check

In this experiment, we examine the legality of personalized license plates by using our developed model and classification system. The legality check focuses on evaluating whether a given license plate meets the established legal and regulatory standards set by various states. This process involves identifying potentially harmful, offensive, or inappropriate content that may violate public decency, safety regulations, or intellectual property laws.

#### 5.2.1 Baseline

The Lexicon-based method for detecting toxic content in personalized license plates combines two features: the presence of toxic words from a predefined list and N-gram features extracted using TfidfVectorizer. These features are merged into a matrix and used to train a Logistic Regression model. The model is evaluated using classification metrics and AUC to assess its ability to classify plates as toxic or non-toxic.

#### 5.2.3 LLM-based Approaches

We conducted experiments using the LLaMA 3.1-8b-instruct model accessed via Ollama[14] and fine-tuned with LoRA through unsloth[15], using 1 epoch and a learning rate of 2e-4. Unsloth provided 2x faster performance and reduced memory usage by 60%. The experiments were performed on an NVIDIA GeForce 4060 GPU and in the Colab GPU-T4 environment.

For the standard prompt design, we structured it into three components: role setting, task description, and output template. In designing the Chain-of-Thought (CoT) prompt, we included examples of step-by-step reasoning to encourage the model to utilize its reasoning capabilities. The process involved having the model first review the license plate content, categorize it, verify that it did not belong to any subcategory of a broader classification, and then provide the appropriate result. The details of the prompt are provided in Appendix B.

5.2.4 Evaluation Metrics

The model's performance was evaluated using standard classification metrics, including accuracy, precision, recall, and F1 score. Additionally, the AUC (Area Under the Curve) was computed to assess the model's ability to distinguish between toxic and non-toxic license plates. The results were compared to the baseline lexicon-based method to highlight improvements in performance.

6. Result

6.1 Identification of Personalized License Plates Result

The pattern-based approach, which uses predefined license plate formats for California and New York, achieved an impressive accuracy of 99%. Additionally, the model maintained high precision, recall, and F1-score, all at 99%, indicating that the method effectively identified personalized plates while minimizing false positives and false negatives. This result confirms that the pattern-based method is highly reliable for distinguishing random license plates from personalized ones.

Model	Result			
	Accuracy	Precision	Recall	F1-Score
Pattern	0.99	0.99	0.99	0.99

Table 2:  
Table 2 presents the performance metrics for the identification of personalized license plates using the pattern-based method. The results demonstrate excellent performance across all evaluation metrics.

Table 2: Identification of Personalized License Plates

6.2 LLM-based Approaches Result

The **Lexicon-based** method, which relies on predefined word lists for detecting illegal terms, performed poorly, with an accuracy of only 49%, and a significantly lower F1-score of 0.43. This suggests that the Lexicon-based method struggles to effectively detect all illegal plates, likely due to its inability to capture the complexity and variety of personalized plate content. As time progresses, the creativity and complexity of personalized license plates have also increased, which further lowers the accuracy of the Lexicon-based approach, as it is unable to account for new and innovative plate combinations.

Among the LLaMA-based models, **LLaMA3 + LoRA** achieved the highest overall accuracy (71%) and balanced performance across all metrics (precision, recall, and F1-score). This demonstrates the effectiveness of combining LLaMA3 with LoRA for fine-tuning, which appears to improve the model's ability to identify legal and illegal personalized plates.

On the other hand, **LLaMA3 + CoT** achieved the highest precision (66%), indicating its strong ability to identify clearly illegal plates. However, its recall and F1-score were somewhat lower, suggesting it may miss some cases. **LLaMA3 + CoT** also achieved an accuracy of 50% on taxonomy. The standard **LLaMA3** model, while not as effective as the fine-tuned versions, still demonstrated reasonable performance, with balanced accuracy, precision, recall, and F1-score all around 59%.

Overall, the results highlight that while the Lexicon-based approach is insufficient, the use of LLaMA3 with fine-tuning techniques such as LoRA or CoT significantly enhances performance, particularly in terms of

accuracy and precision. The combination of LLaMA3 + LoRA stands out as the most well-rounded approach for identifying the legality of personalized plates.

Model	Result			
	Accuracy	Precision	Recall	F1-Score
Lexicon-based	0.49	0.48	0.49	0.43
llama3	0.59	0.59	0.59	0.59
llama3+CoT	0.60	0.66	0.60	0.56
llama3+LoRA	<b>0.71</b>	<b>0.72</b>	<b>0.69</b>	<b>0.70</b>

Table 3: LLM-based Approaches

Table 3:  
This presents the performance metrics for identifying whether personalized license plates are legal using various LLM-based approaches. The results indicate that different methods show varying levels of effectiveness in terms of accuracy, precision, recall, and F1-score.

## 7. Discussion

Although LLaMA3 + LoRA achieved the highest performance in our experiments with an accuracy of 71%, its overall effectiveness is still relatively modest. There are several possible reasons for this performance.

1. Limited Model Size  
LLaMA3.1, with only 7 billion parameters, may be too small for the complexity of the task at hand. While it performs well in many natural language processing tasks, personalized license plate classification, especially in the context of identifying illegal plates, may require a model with more parameters to better capture the nuances and variations of plate content. A larger model could potentially lead to better generalization and more accurate results.
2. Complexity of License Plate Abbreviations  
Personalized license plates often contain abbreviations, slang, or culturally specific references that can be challenging for models to interpret. Even with fine-tuning, LLaMA3 might still struggle to consistently understand and classify such intricate content, particularly without a strong contextual understanding or explicit knowledge base for slang and abbreviations.
3. Data Quality and Imbalance  
Although efforts were made to balance the dataset, there might still be issues with data quality and representativeness. For example, in the California dataset, reviewer comments indicate that some plates containing terms like "red" or "gang color" were accepted, while others with similar content were rejected. This inconsistency may reflect subjective judgment by different reviewers, introducing noise and variability into the dataset. This variability in annotations could significantly affect model training, as the model may learn conflicting signals from the data.

In conclusion, the combination of LLaMA3 + LoRA shows potential but is hindered by several limitations. The model's relatively small size, the complexity of interpreting personalized license plate content, and inconsistencies in the dataset all contribute to the modest performance. Despite this, the approach demonstrates that LLaMA3 can be a useful tool in the classification and detection of personalized license plates, though further improvements, such as larger models, more domain-specific fine-tuning, and cleaner data, are necessary for achieving better accuracy and reliability. Future work could explore alternative approaches to further enhance model performance, including the use of more advanced fine-tuning techniques or the incorporation of additional domain-specific knowledge.



## Reference

- [1] New York State Department of Motor Vehicles, "Restrictions on personalized plates," DMV, [Online]. Available: <https://dmv.ny.gov/plates/custom-plates/restrictions-on-personalized-plates>. [Accessed: Dec. 5, 2024].
- [2] R. Pod, "The psychology of vanity plates: What your license plate says about you," Rusty Pod, [Online]. Available: <https://rustypod.com/the-psychology-of-vanity-plates-what-your-license-plate-says-about-you/>. [Accessed: Dec. 5, 2024].
- [3] H. Inan, K. Upasani, J. Chi, R. Rungta, K. Iyer, Y. Mao, M. Tontchev, Q. Hu, B. Fuller, D. Testuggine, and M. Khabsa, "Llama Guard: LLM-based input-output safeguard for human-AI conversations," arXiv, Dec. 7, 2023. [Online]. Available: <https://arxiv.org/abs/2312.06674>. [Accessed: Dec. 5, 2024].
- [4] California Department of Motor Vehicles, "Personalized configurations mandatory refusal," DMV California, [Online]. Available: <https://www.dmv.ca.gov/portal/handbook/vehicle-industry-registration-procedures-manual-2/special-plates/personalized-configurations-mandatory-refusal/>. [Accessed: Dec. 5, 2024].
- [5] Veltman, "ca-license-plates," GitHub, [Online]. Available: <https://github.com/veltman/ca-license-plates>. [Accessed: Dec. 5, 2024].
- [6] Datanews, "license-plates," GitHub, [Online]. Available: <https://github.com/datanews/license-plates>. [Accessed: Dec. 5, 2024].
- [7] Orthrus-Lexicon, "Toxic," GitHub, [Online]. Available: <https://github.com/Orthrus-Lexicon/Toxic>. [Accessed: Dec. 5, 2024].
- [8] Y. Xiao, Y. Hu, K. T. W. Choo, and R. K.-w. Lee, "ToxiCloakCN: Evaluating robustness of offensive language detection in Chinese with cloaking perturbations," arXiv, Jun. 18, 2024. [Online]. Available: <https://arxiv.org/abs/2406>.
- [9] E. J. Hu, Y. Shen, P. Wallis, Z. Allen-Zhu, Y. Li, S. Wang, L. Wang, and W. Chen, "Lora: Low-rank adaptation of large language models," arXiv preprint arXiv:2106.09685, 2021.
- [10] J. Wei, X. Wang, D. Schuurmans, M. Bosma, B. Ichter, F. Xia, E. Chi, Q. Le, and D. Zhou, "Chain-of-thought prompting elicits reasoning in large language models," arXiv, Jan. 28, 2022. [Online]. Available: <https://arxiv.org/abs/2201.11903>. [Accessed: Dec. 5, 2024].
- [11] T. Sen, A. Das, and M. Sen, "HateTinyLLM: Hate speech detection using tiny large language models," arXiv, Apr. 26, 2024. [Online]. Available: <https://arxiv.org/abs/2405.01577>. [Accessed: Dec. 5, 2024].
- [12] R. O. D. K. Stitt, "Study: 83% of accidents result from drivers deciphering vanity license plates," The Hard Times, Dec. 4, 2024. [Online]. Available: <https://thehardtimes.net/culture/study-83-of-accidents-result-from-drivers-deciphering-vanity-license-plates/>. [Accessed: Dec. 5, 2024].

- There is a high frequency of words related to the owner and vehicle's attributes, such as "name," "nickname," "year," "model," "initial," and "last."

- There are also words expressing the owner's personal preferences, such as color choices and lucky numbers.
- Additionally, some words are related to the owner's profession or company names.

## B. Prompts

### Standard prompting

You are an officer in the Department of Motor Vehicles reviewing a personalized license plate {plate}. Based on content and appropriateness, should the license plate be approved or denied? The format of the license plate is not a consideration. You need to consider all the tricks the applicant may use to conceal his or her intentions.

Please respond with a JSON object containing the result ("Approve", "Disapprove") and a brief explanation. If the response contains any inappropriate words, replace them with asterisks (\*):

"result": "Approve/Disapprove",

"explanation": "Brief explanation of the decision."

### Chain-of-thought prompting

You are a DMV officer tasked with reviewing personalized license plate requests. For each license plate, you need to categorize it as "Approve" or "Disapprove" based on the following categories. Your final decision should be provided in a JSON format containing the "result," "category," and "reason" fields.

**Rejected License Plates:**

1. Contains illegal or violent content.
2. Contains sexual or inappropriate content.
3. Contains insulting, hateful, or discriminatory content.
4. Contains intellectual property violations or could cause confusion.
5. Restricted due to technical or special-use limitations.
6. Others

**Approved License Plates:**

1. Personal identification (e.g., names, nicknames, professions).
2. Vehicle-related (e.g., color, model, year).
3. Hobbies and interests.
4. Creative or humorous expressions.
5. Profession or organization-related (e.g., job titles, company names).
6. Others

**Example:**

You receive a license plate request for "SMOKER1".

1. Review the content: The license plate "SMOKER1" could imply a smoker or someone who smokes.
2. Check the categories: This could be linked to personal identification (e.g., referring to a personal habit).
3. Assess if it fits any rejection category: This does not contain illegal, violent, sexual, insulting, hateful, or discriminatory content. It is not related to intellectual property or technical restrictions.
4. Final Decision (JSON format):

"result": "Approve",

"category": 1,

"reason": "The license plate relates to personal identification or hobbies/interests."

Now, you will need to follow the same process for each new license plate request. After reviewing the plate and applying the categories, provide your final decision in the same JSON format, with a brief explanation.

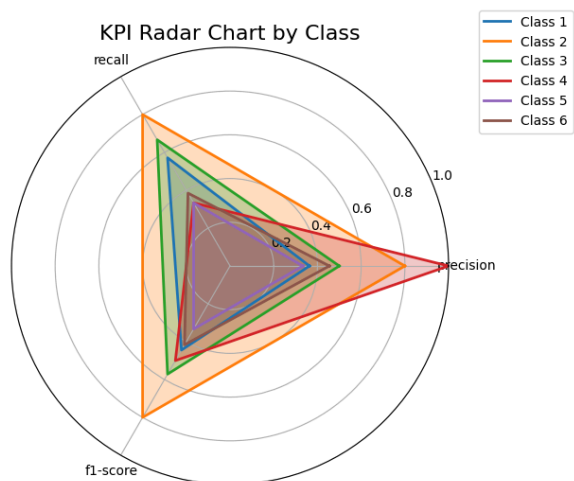
**Input:**

- License plate text: {plate}

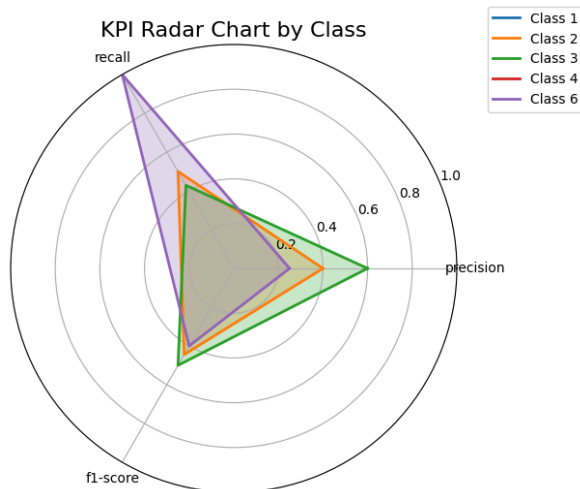
**Output:**

### C. Categories

The classification accuracy for the categories was observed to be 0.5, indicating that the model correctly identified half of the cases. When broken down by true positive and true negative rates, the results highlight an imbalance in performance. The true positive accuracy reached 0.53, showing the model's relatively better ability to correctly identify instances belonging to positive categories. However, the true negative accuracy was only 0.4, suggesting a notable difficulty in accurately classifying negative cases.



True Positive KPI Plot



True Negative KPI plot

The KPI charts show that most categories have balanced R, P, and F1 scores. However, Class 4 (intellectual property or misleading content) for approved plates has exceptionally high precision, indicating the model is highly cautious in identifying this category with minimal false positives. Conversely, Class 6 (others) for rejected plates has very high recall, suggesting the model captures nearly all instances of this class but may include more false positives, reducing precision.