






On the Cross-Dataset Generalization in License Plate Recognition

Rayson Laroca¹, Everton V. Cardoso¹, Diego R. Lucio¹,
Valter Estevam^{1,2}, and David Menotti¹

¹Federal University of Paraná, Curitiba, Brazil

²Federal Institute of Paraná, Irati, Brazil

{rblsantos, ecardoso, drlucio, vlejunior, menotti}@inf.ufpr.br, valter.junior@ifpr.edu.br

Keywords: Deep Learning, Leave-one-dataset-out, License Plate Recognition, Optical Character Recognition

Abstract: Automatic License Plate Recognition (ALPR) systems have shown remarkable performance on license plates (LPs) from multiple regions due to advances in deep learning and the increasing availability of datasets. The evaluation of deep ALPR systems is usually done within each dataset; therefore, it is questionable if such results are a reliable indicator of generalization ability. In this paper, we propose a traditional-split *versus* leave-one-dataset-out experimental setup to empirically assess the cross-dataset generalization of 12 Optical Character Recognition (OCR) models applied to LP recognition on nine publicly available datasets with a great variety in several aspects (e.g., acquisition settings, image resolution, and LP layouts). We also introduce a public dataset for end-to-end ALPR that is the first to contain images of vehicles with Mercosur LPs and the one with the highest number of motorcycle images. The experimental results shed light on the limitations of the traditional-split protocol for evaluating approaches in the ALPR context, as there are significant drops in performance for most datasets when training and testing the models in a leave-one-dataset-out fashion.

1 INTRODUCTION

The global automotive industry expects to produce more than 82 million light vehicles in 2022 alone, despite the ongoing coronavirus pandemic and chip supply issues (Forbes, 2021; IHS Markit, 2021). In addition to bringing convenience to owners, vehicles also significantly modify the urban environment, posing challenges concerning pollution, privacy and security – especially in large urban centers. The constant monitoring of vehicles through computational techniques is of paramount importance and, therefore, it has been a frequent research topic. In this context, Automatic License Plate Recognition (ALPR) systems (Weihong and Jiaoyang, 2020; Lubna et al., 2021) stand out.

ALPR systems exploit image processing and pattern recognition techniques to detect and recognize the characters on license plates (LPs) from images or videos. Some practical applications for an ALPR system are road traffic monitoring, toll collection, and vehicle access control in restricted areas (Špaňhel et al.,

2017; Henry et al., 2020; Wang et al., 2022).

Deep ALPR systems have shown remarkable performance on LPs from multiple regions due to advances in deep learning and the increasing availability of datasets (Henry et al., 2020; Silva and Jung, 2022). In the past, the evaluation of ALPR systems used to be done within each of the chosen datasets, i.e., the proposed methods were trained and evaluated on different subsets from the same dataset. Such an evaluation was carried out independently for each dataset. Recently, considering that deep models can take considerable time to be trained (especially on low- or mid-end GPUs), the authors have adopted a protocol where the proposed models are trained once on the union of the training images from the chosen datasets and evaluated individually on the respective test sets (Selmi et al., 2020; Laroca et al., 2021b). Although the images for training and testing belong to disjoint subsets, these protocols do not make it clear whether the evaluated models have good generalization ability, i.e., whether they perform well on images from other scenarios, mainly due to domain divergence and data selection bias (Torralba and Efros, 2011; Tommasi et al., 2017; Zhang et al., 2019).

In this regard, many computer vision researchers

This is an author-prepared version of a paper accepted for presentation at the International Conference on Computer Vision Theory and Applications (VISAPP) 2022. The published version is available at the *SciTePress Digital Library* (DOI: [10.5220/0010846800003124](https://doi.org/10.5220/0010846800003124)).

have carried out cross-dataset experiments – where training and testing data come from different sources – to assess whether the proposed models perform well on data from an unknown domain (Ashraf et al., 2018; Zhang et al., 2019; Estevam et al., 2021). However, as far as we know, there is no work focused on such experimental settings in the ALPR context.

Considering the above discussion, in this work we evaluate for the first time various Optical Character Recognition (OCR) models for LP recognition in a leave-one-dataset-out experimental setup over nine public datasets with different characteristics. The results obtained are compared with those achieved when training the models in the same way as in recent works, that is, using the union of the training set images from all datasets (hereinafter, this protocol is referred to as traditional-split).

Deep learning-based ALPR systems have often achieved recognition rates above 99% in existing datasets under the traditional-split protocol (some examples are provided in Section 2). However, in real-world applications, new cameras are regularly being installed in new locations without existing systems being retrained as often, which can dramatically decrease the performance of those models. A leave-one-dataset-out protocol enables simulating this specific scenario and providing an adequate evaluation of the generalizability of the models.

ALPR is commonly divided into two tasks: LP detection and LP recognition. The former refers to locating the LP region in the input image, while the latter refers to extracting the string related to the LP. In this work, we focus on the LP recognition stage since it is the current bottleneck of ALPR systems (Laroca et al., 2021b). Thus, we simply train the off-the-shelf YOLOv4 model (Bochkovskiy et al., 2020) to detect the LPs in the input images. For completeness, we also report the results achieved in this stage on both of the aforementioned protocols.

As part of this work, we introduce a publicly available dataset, called RodoSol-ALPR¹, that contains 20,000 images captured at toll booths installed on a Brazilian highway. It has images of two different LP layouts: Brazilian and Mercosur², with half of the vehicles being motorcycles (see details in Section 3). To the best of our knowledge, this is the first public dataset for ALPR with images of Mercosur LPs and the largest in the number of motorcycle images. This

¹The RodoSol-ALPR dataset is publicly available to the research community at <https://github.com/raysonlaroca/rodosol-alpr-dataset/>

²Mercosur (*Mercado Común del Sur*, i.e., Southern Common Market in Castilian) is an economic and political bloc comprising Argentina, Brazil, Paraguay and Uruguay.

last information is relevant because motorcycle LPs have two rows of characters, which is a challenge for sequential/recurrent-based methods (Silva and Jung, 2022), and therefore have been overlooked in the evaluation of LP recognition models (see Section 2).

Our paper has two main contributions:

- A traditional-split *versus* leave-one-dataset-out experimental setup that can be considered a valid testbed for cross-dataset generalization methods proposed in future works on ALPR. We present a comparative assessment of 12 OCR models for LP recognition on nine publicly available datasets. The main findings were that (i) there are significant drops in performance for most datasets when training and testing the recognition models in a leave-one-dataset-out fashion, especially when there are different fonts of characters in the training and test images; (ii) no model achieved the best result in all experiments, with 6 different models reaching the best result in at least one dataset under the leave-one-dataset-out protocol; and (iii) the proposed dataset proved very challenging, as both the models trained by us and two commercial systems failed to reach recognition rates above 70% on its test set images.
- A public dataset with 20,000 images acquired in real-world scenarios, being half of them of vehicles with Mercosur LPs. Indeed, one of the objectives of this work is to provide a reliable source of information about Mercosur LPs, as much news – often outdated – has been used as references.

The remainder of this paper is organized in the following manner. In Section 2, we briefly review related works. The RodoSol-ALPR dataset is introduced in Section 3. The setup adopted in our experiments is thoroughly described in Section 4. Section 5 presents the results achieved. Finally, Section 6 concludes the paper and outlines future directions of research.

2 RELATED WORK

In this section, we first present concisely recent works on LP recognition. Then, we situate the current state of ALPR research in terms of cross-dataset experiments, Mercosur LPs, and motorcycle LPs.

The good speed/accuracy trade-off provided by YOLO networks (Redmon et al., 2016; Bochkovskiy et al., 2020) inspired many authors to explore similar architectures targeting real-time performance for LP recognition. For example, Silva and Jung (2020) proposed a YOLO-based model to simultaneously detect and recognize all characters within a cropped LP. This

model, called CR-NET, consists of the first eleven layers of YOLO and four other convolutional layers added to improve non-linearity. Impressive results were achieved through CR-NET both in the original work and in more recent ones (Laroca et al., 2021b; Oliveira et al., 2021; Silva and Jung, 2022).

While Kessentini et al. (2019) applied the YOLOv2 model without any change or refinement to this task, Henry et al. (2020) used a modified version of YOLOv3 that includes spatial pyramid pooling. Although these two models achieved high recognition rates in multiple datasets, they are very deep for LPs recognition, making it difficult to meet the real-time requirements of ALPR applications.

Rather than exploring object detectors, Zou et al. (2020) adopted a bi-directional Long Short-Term Memory (LSTM) network to implicitly locate the characters on the LP. They explored a 1-D attention module to extract useful features of the character regions, improving the accuracy of LP recognition. In a similar way, Zhang et al. (2021) used a 2-D attention mechanism to optimize their recognition model, which uses a 30-layer Convolutional Neural Network (CNN) based on Xception for feature extraction. An LSTM model was adopted to decode the extracted features into LP characters.

There are also several works where multi-task networks were designed to holistically process the entire LP image and, thus, avoid character segmentation, such as (Špaňhel et al., 2017; Gonçalves et al., 2019). As these networks employ fully connected layers as classifiers to recognize the characters on the predefined positions of the LPs, they may not generalize well with small-scale training sets since the probability of a specific character appearing in a specific position is low. To deal with this, Wang et al. (2022) proposed a weight-sharing classifier, which is able to spot instances of each character across all positions.

Considering that the recognition rates achieved under the traditional-split protocol have significantly increased in recent years, some authors began to conduct small cross-dataset experiments to analyze the generalization ability of the proposed methods. For example, Silva and Jung (2020); Laroca et al. (2021b) used all 108 images from the OpenALPR-EU dataset for testing, rather than using some for training/validation. Nevertheless, the results achieved in so few test images are susceptible to tricks, especially considering that heuristic rules were explored to improve the LP recognition results in both works.

As another example, Zou et al. (2020); Zhang et al. (2021); Wang et al. (2022) trained their recognition models specifically for Chinese LPs on approximately 200K images from the CCPD dataset (Xu

et al., 2018) and tested them on images from other datasets that also contain only Chinese LPs. In this case, it is not clear whether the proposed models perform well on LPs from other regions. In fact, the authors trained another instance of the respective models to evaluate them in the AOLP dataset (Hsu et al., 2013), which contains LPs from the Taiwan region.

Recently, Mercosur countries adopted a unified standard of LPs for newly purchased vehicles, inspired by the integrated system adopted by European Union countries many years ago. Although the new standard has been implemented in all countries in the bloc, there is still no public dataset for ALPR with images of Mercosur LPs as far as we know.

In this sense, Silvano et al. (2021) presented a methodology that couples synthetic images of Mercosur LPs with real-world images containing vehicles with other LP layouts. A model trained exclusively with synthetic images achieved promising results on 1,000 real images from various sources; however, it is difficult to assess these results accurately since the test images were not made available to the research community. The LP recognition stage was not addressed.

Despite the fact that motorcycles are one of the most popular transportation means in metropolitan areas (Hsu et al., 2015), they have been largely overlooked in ALPR research. There are even works where images of motorcycles were excluded from the experiments (Gonçalves et al., 2018; Silva and Jung, 2020), mainly because LPs of motorcycles usually have two rows of characters, which are challenging to sequential/recurrent-based methods (Kessentini et al., 2019; Silva and Jung, 2022), and also because they are generally smaller in size (having less space between characters) and are often tilted.

In this regard, there is a great demand for a public dataset for end-to-end ALPR with the same number of images of cars and motorcycles to give equal importance to LPs with one or two rows of characters in the assessment of ALPR systems.

3 RODOSOL-ALPR DATASET

The RodoSol-ALPR dataset contains 20,000 images captured by static cameras located at pay tolls owned by the *Rodovia do Sol* (RodoSol) concessionaire (RodoSol, 2022) (hence the name of the dataset), which operates 67.5 kilometers of a highway (ES-060) in the Brazilian state of Espírito Santo.

As can be seen in Figure 1, there are images of different types of vehicles (e.g., cars, motorcycles, buses and trucks), captured during the day and night, from distinct lanes, on clear and rainy days, and the dis-



Figure 1: Some images extracted from the RodoSol-ALPR dataset. The first and second rows show images of cars and motorcycles, respectively, with Brazilian LPs (i.e., the standard used in Brazil before the adoption of the Mercosur standard). The third and fourth rows show images of cars and motorcycles, respectively, with Mercosur LPs. We show a zoomed-in version of the vehicle’s LP in the lower right region of the images in the last column for better viewing of the LP layouts.

tance from the vehicle to the camera varies slightly. All images have a resolution of $1,280 \times 720$ pixels.

An important feature of the proposed dataset is that it has images of two different LP layouts: Brazilian and Mercosur. To maintain consistency with previous works (Izidio et al., 2020; Oliveira et al., 2021; Silva and Jung, 2022), we refer to “Brazilian” as the standard used in Brazil before the adoption of the Mercosur standard. All Brazilian LPs consist of three letters followed by four digits, while the initial pattern adopted in Brazil for Mercosur LPs consists of three letters, one digit, one letter and two digits, in that order. In both layouts, car LPs have seven characters arranged in one row, whereas motorcycle LPs have three characters in one row and four characters in another. Even though these LP layouts are very similar in shape and size, there are considerable differences in their colors and characters’ fonts.

The 20,000 images are divided as follows: 5,000 images of cars with Brazilian LPs; 5,000 images of motorcycles with Brazilian LPs; 5,000 images of cars with Mercosur LPs; and 5,000 images of motorcycles with Mercosur LPs. For the sake of simplicity of definitions, here “car” refers to any vehicle with four wheels or more (e.g., passenger cars, vans, buses, trucks, among others), while “motorcycle” refers to both motorcycles and motorized tricycles. As far as we know, RodoSol-ALPR is the public dataset for ALPR with the highest number of motorcycle images.

We randomly split the RodoSol-ALPR dataset as follows: 8,000 images for training; 8,000 images for testing; and 4,000 images for validation, following the split protocol (i.e., 40%/40%/20%) adopted in the

SSIG-SegPlate (Gonçalves et al., 2016) and UFPR-ALPR (Laroca et al., 2018) datasets. We preserved the percentage of samples for each vehicle type and LP layout; for example, there are 2,000 images of cars with Brazilian LPs in each of the training and test sets, and 1,000 images in the validation one. For reproducibility purposes, the subsets generated are explicitly available along with the proposed dataset.

Every image has the following information available in a text file: the vehicle’s type (car or motorcycle), the LP’s layout (Brazilian or Mercosur), its text (e.g., ABC-1234), and the position (x, y) of each of its four corners. We labeled the corners instead of just the LP bounding box to enable the training of methods that explore LP rectification, as well as the application of a wider range of data augmentation techniques.

The datasets for ALPR are generally very unbalanced in terms of character classes due to LP allocation policies (Zhang et al., 2021). In Brazil, for example, one letter can appear much more often than others according to the state in which the LP was issued (Gonçalves et al., 2018; Laroca et al., 2018). This information must be taken into account when training recognition models in order to avoid undesirable biases – this is usually done through data augmentation techniques (Zhang et al., 2021; Hasnat and Nakib, 2021); for example, a network trained exclusively in our dataset may learn to always classify the first character as ‘P’ in cases where it should be ‘B’ or ‘R’ since it appears much more often in this position than these two characters (see Figure 2).

Regarding privacy concerns related to our dataset, we remark that in Brazil the LPs are related to the re-

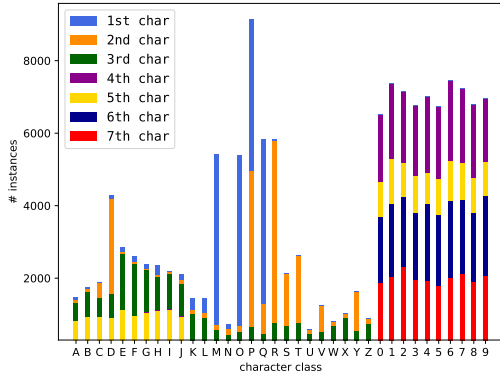


Figure 2: The distribution of character classes in the RodoSol-ALPR dataset. Observe that there is a significant imbalance in the distribution of the letters (due to LP allocation policies), whereas the digits are well balanced.

spective vehicles, i.e., no public information is available about the vehicle drivers/owners (Presidência da República, 1997; Oliveira et al., 2021). Moreover, all human faces (e.g., drivers or RodoSol’s employees) were manually redacted (i.e., blurred) in each image.

4 EXPERIMENTS

In this section, we describe the setup adopted in our experiments. We first list the models we implemented for our assessments, explaining why they were chosen and not others. Afterward, we provide the implementation details, for example, which framework was used to train/test each model and the respective hyperparameters. We then present and briefly describe the datasets used in our experiments, as well as the data augmentation techniques explored to avoid overfitting. Lastly, we detail the evaluation protocols adopted by us, that is, which images from each dataset were used for training or testing in each experiment, and how we evaluate the performance of each method.

4.1 Methods

In this work, we evaluate 12 OCR models for LP recognition: RARE (Shi et al., 2016), R²AM (Lee and Osindero, 2016), STAR-Net (Liu et al., 2016), CRNN (Shi et al., 2017), GRCNN (Wang and Hu, 2017), Holistic-CNN (Špaňhel et al., 2017), Multi-task (Gonçalves et al., 2019), Rosetta (Borisjuk et al., 2018), TRBA (Baek et al., 2019), CR-NET (Silva and Jung, 2020), Fast-OCR (Laroca et al., 2021a), and ViTSTR-Base (Atienza, 2021). Table 1 presents an overview of these methods, listing the original OCR application for which they were designed as well as the framework we used to train and evaluate them.

Table 1: OCR models explored in our experiments.

| Model | Original Application |
|--|---------------------------|
| Framework: PyTorch ³ | |
| R ² AM (Lee and Osindero, 2016) | Scene Text Recognition |
| RARE (Shi et al., 2016) | Scene Text Recognition |
| STAR-Net (Liu et al., 2016) | Scene Text Recognition |
| CRNN (Shi et al., 2017) | Scene Text Recognition |
| GRCNN (Wang and Hu, 2017) | Scene Text Recognition |
| Rosetta (Borisjuk et al., 2018) | Scene Text Recognition |
| TRBA (Baek et al., 2019) | Scene Text Recognition |
| ViTSTR-Base (Atienza, 2021) | Scene Text Recognition |
| Framework: Keras ⁴ | |
| Holistic-CNN (Špaňhel et al., 2017) | License Plate Recognition |
| Multi-task (Gonçalves et al., 2019) | License Plate Recognition |
| Framework: Darknet ⁵ | |
| CR-NET (Silva and Jung, 2020) | License Plate Recognition |
| Fast-OCR (Laroca et al., 2021a) | Image-based Meter Reading |

These models were chosen/implemented by us for two main reasons: (i) they have been employed for OCR tasks with promising/impressive results (Baek et al., 2019; Atienza, 2021; Laroca et al., 2021a), and (ii) we believe we have the necessary knowledge to train/adjust them in the best possible way in order to ensure fairness in our experiments, as the authors provided enough details about the architectures used, and also because we designed/employed similar networks in previous works (even the same ones in some cases) (Gonçalves et al., 2018, 2019; Laroca et al., 2019, 2021a). Note that we are not aware of any work in the ALPR literature where so many recognition models were explored in the experiments.

The CR-NET and Fast-OCR models are based on the YOLO object detector (Redmon et al., 2016). Thus, they are trained to predict 35 classes (0-9, A-Z, where ‘O’ and ‘o’ are detected/recognized jointly) using the bounding box of each LP character as input. Although these methods have been attaining impressive results, they require laborious data annotations, i.e., each character’s bounding box needs to be labeled for training them (Wang et al., 2022). All the other 10 models, on the other hand, output the LP characters in a segmentation-free manner, i.e., they predict the characters holistically from the LP region without the need to detect/segment them. According to previous works (Gonçalves et al., 2018; Atienza, 2021; Hasnat and Nakib, 2021), the generalizability of such segmentation-free models tends to improve significantly through the use of data augmentation.

³<https://github.com/roatienza/deep-text-recognition-benchmark/>

⁴<https://keras.io/>

⁵<https://github.com/AlexeyAB/darknet/>

4.2 Setup

All experiments were carried out on a computer with an AMD Ryzen Threadripper 1920X 3.5GHz CPU, 96 GB of RAM (2133 MHz), HDD 7200 RPM, and an NVIDIA Quadro RTX 8000 GPU (48 GB).

Although run-time analysis is considered a critical factor in the ALPR literature (Lubna et al., 2021), we consider such analysis beyond the scope of this work since we used different frameworks to implement the recognition models and there are probably differences in implementation and optimization between them – we implemented each method using either the framework where it was originally implemented or well-known public repositories. For example, the YOLO-based models were implemented using Darknet⁵ while the models originally proposed for scene text recognition were trained and evaluated using a fork³ of the open source repository of Clova AI Research (PyTorch) used to record the 1st place of ICDAR2013 focused scene text and ICDAR2019 ArT, and 3rd place of ICDAR2017 COCO-Text and ICDAR2019 ReCTS (task1) (Baek et al., 2019).

For completeness, below we list the hyperparameters used in each framework for training the OCR models; we remark that these hyperparameters were defined based on previous works as well as on experiments performed in the validation set. In Darknet, we employed the following parameters: Stochastic Gradient Descent (SGD) optimizer, 90K iterations (max batches), batch size = 64, and learning rate = $[10^{-3}, 10^{-4}, 10^{-5}]$ with decay steps at 30K and 60K iterations. In Keras, we used the Adam optimizer, initial learning rate = 10^{-3} (with *ReduceLROnPlateau*’s patience = 5 and factor = 10^{-1}), batch size = 64, max epochs = 100, and patience = 11 (patience refers to the number of epochs with no improvement after which training is stopped). In PyTorch, we adopted the following parameters: Adadelta optimizer, whose decay rate is set to $\rho = 0.99$, 300K iterations, and batch size = 128.

4.3 Datasets

Our experiments were conducted on images from the RodoSol-ALPR dataset and eight publicly available datasets that are often employed to benchmark ALPR algorithms: Caltech Cars (Weber, 1999), EnglishLP (Srebrić, 2003), UCSD-Stills (Dlagnakov and Belongie, 2005), ChineseLP (Zhou et al., 2012), AOLP (Hsu et al., 2013), OpenALPR-EU (OpenALPR Inc., 2016), SSIG-SegPlate (Gonçalves et al., 2016), UFPR-ALPR (Laroca et al., 2018). Table 2 shows an overview of these datasets. They

were introduced over the last 23 years and have considerable diversity in terms of the number of images, acquisition settings, image resolution, and LP layouts. As far as we know, there is no other work in the ALPR literature where experiments were carried out on images from so many public datasets.

Table 2: The datasets used in our experiments. In this work, the “Chinese” layout refers to LPs of vehicles registered in mainland China, while the “Taiwanese” layout refers to LPs of vehicles registered in the Taiwan region.

| Dataset | Year | Images | Resolution | LP Layout |
|---------------|------|--------|--------------------|--------------------|
| Caltech Cars | 1999 | 126 | 896×592 | American |
| EnglishLP | 2003 | 509 | 640×480 | European |
| UCSD-Stills | 2005 | 291 | 640×480 | American |
| ChineseLP | 2012 | 411 | Various | Chinese |
| AOLP | 2013 | 2049 | Various | Taiwanese |
| OpenALPR-EU | 2016 | 108 | Various | European |
| SSIG-SegPlate | 2016 | 2000 | 1920×1080 | Brazilian |
| UFPR-ALPR | 2018 | 4500 | 1920×1080 | Brazilian |
| RodoSol-ALPR | 2022 | 20000 | 1280×720 | Brazilian/Mercosur |

Figure 3 shows the diversity of the chosen datasets in terms of LP layouts. It is clear that even LPs from the same country can be quite different, e.g., the Caltech Cars and UCSD-Stills datasets were collected in the same region (California, United States), but they have images of LPs with significant differences in terms of colors, aspect ratios, backgrounds, and the number of characters. It can also be observed that some datasets have LPs with two rows of characters and that the LPs may be tilted or have low resolution due to camera quality or vehicle-to-camera distance.

In order to eliminate biases from the public datasets, we also used 772 images from the internet – those labeled and provided by Laroca et al. (2021b) – to train all models. These images include 257 American LPs, 347 Chinese LPs, and 178 European LPs. We chose not to use two datasets introduced recently: KarPlate (Henry et al., 2020) and CCPD (Xu et al., 2018). The former cannot currently be downloaded due to legal problems. The latter, although already available, was not employed for two main reasons: (i) it contains highly compressed images, which significantly compromises the readability of the LPs (Silva and Jung, 2022); and (ii) it has some large errors in the corners’ annotations (Meng et al., 2020) – this was somewhat expected since the corners were labeled automatically using RPnet (Xu et al., 2018). Additionally, we could not download the CLPD dataset (Zhang et al., 2021), as the authors made it available exclusively through a Chinese website where registration – using a Chinese phone number or identity document – is required (we contacted the authors requesting an alternative link to download the dataset, but have not received a response so far).



Figure 3: Some LP images from the public datasets used in our experimental evaluation. We show some LP images from the RodoSol-ALPR dataset in the last column of Fig 1.

4.3.1 Data Augmentation

As shown in Table 2, two-thirds of the images used in our experiments are from the RodoSol-ALPR dataset. In order to prevent overfitting, we initially balanced the number of images from different datasets through data augmentation techniques such as random cropping, random shadows, conversion to grayscale, and random perturbations of hue, saturation and brightness. We used Albumentations (Buslaev et al., 2020), which is a well-known Python library for image augmentation, to apply these transformations. Nevertheless, preliminary experiments showed that some of the recognition models were prone to predict only LP patterns that existed in the training set, as some patterns were being fed numerous times per epoch to the networks – especially from small-scale datasets, where many images were created from a single original one. Therefore, inspired by Gonçalves et al. (2018), we also randomly permuted the position of the characters on each LP to eliminate such biases in the learning process (as illustrated in Figure 4). As the bounding box of each LP character is required to apply this data augmentation technique – these annotations are very time-consuming and laborious – we do not augment the training images from the RodoSol-ALPR dataset. We believe this is not a significant problem as the proposed dataset is much larger than the others. The images from the other public datasets were augmented

using the labels provided by Laroca et al. (2021b).



Figure 4: Illustration of the character permutation-based data augmentation technique (Gonçalves et al., 2018) we adopted to avoid overfitting. The images in the first row are the originals, while the others were generated automatically.

In this process, we do not enforce the generated LPs to have the same arrangement of letters and digits of the original LPs so that the recognition models do not memorize specific patterns from different LP layouts. For example, as described in Section 3, all Brazilian LPs consist of 3 letters followed by 4 digits, while Mercosur LPs have 3 letters, 1 digit, 1 letter and 2 digits, in that order. Considering that LPs of these layouts are relatively similar (in size, shape, etc.), the segmentation-free networks would probably predict 3 letters followed by 4 digits for most Mercosur LP when holding the RodoSol-ALPR dataset out in a leave-one-dataset-out evaluation, as none of the other datasets have vehicles with Mercosur LPs.

4.4 Evaluation Protocols

In our experiments, we consider both traditional-split and leave-one-dataset-out protocols. In the following subsections, we first describe them in detail. Then, we discuss how the performance evaluation is carried out.

4.4.1 Traditional Split

The traditional-split protocol assesses the ability of the models to perform well in seen scenarios, as each model is trained on the union of the training set images from all datasets and evaluated on the test set images from the respective datasets. In recent works, the authors have chosen to train a single model on images from multiple datasets (instead of training a specific network for each dataset or LP layout as was commonly done in the past) so that the proposed models are robust for different scenarios with considerably less manual effort since their parameters are adjusted only once for all datasets (Selmi et al., 2020; Laroca et al., 2021b; Silva and Jung, 2022).

For reproducibility, it is important to make clear how we divided the images from each of the datasets to train, validate and test the chosen models. The

UCSD-Stills, SSIG-SegPlate, UFPR-ALPR and RodoSol-ALPR datasets were split according to the protocols defined by the respective authors, while the other datasets, which do not have well-defined evaluation protocols, were divided following previous works. In summary, as in (Xiang et al., 2019; Henry et al., 2020), the Caltech Cars dataset was randomly split into 80 images for training/validation and 46 images for testing. Following (Panahi and Gholampour, 2017; Beratoğlu and Töreyn, 2021), the EnglishLP dataset was randomly divided as follows: 80% of the images for training/validation and 20% for testing. For the ChineseLP dataset, we employed the same protocol as Laroca et al. (2021b): 40% of the images for training, 20% for validation and 40% for testing. We split each of the three subsets of the AOLP dataset (i.e., AC, LE, and RP) into training and test sets with a 2:1 ratio, following (Xie et al., 2018; Liang et al., 2021), with 20% of the training images being used for validation. Finally, as most works in the literature (Masood et al., 2017; Laroca et al., 2021b; Silva and Jung, 2022), we used all the 108 images from the OpenALPR-EU dataset for testing (this division has been considered as a mini leave-one-dataset-out evaluation in recent works). Table 3 lists the exact number of images used for training, validating and testing the chosen models.

Table 3: An overview of the number of images from each dataset used for training, validation, and testing.

| Dataset | Training | Validation | Testing | Discarded | Total |
|---------------|----------|------------|---------|-----------|--------|
| Caltech Cars | 61 | 16 | 46 | 3 | 126 |
| EnglishLP | 326 | 81 | 102 | 0 | 509 |
| UCSD-Stills | 181 | 39 | 60 | 11 | 291 |
| ChineseLP | 159 | 79 | 159 | 14 | 411 |
| AOLP | 1,093 | 273 | 683 | 0 | 2,049 |
| OpenALPR-EU | 0 | 0 | 108 | 0 | 108 |
| SSIG-SegPlate | 789 | 407 | 804 | 0 | 2,000 |
| UFPR-ALPR | 1,800 | 900 | 1,800 | 0 | 4,500 |
| RodoSol-ALPR | 8,000 | 4,000 | 8,000 | 0 | 20,000 |

As also detailed in Table 3, a few images (0.01%) were discarded in our experiments because it is impossible to recognize the LP(s) on them due to occlusion, lighting or image acquisition problems⁶. Such images were also discarded by Masood et al. (2017) and Laroca et al. (2021b).

4.4.2 Leave-one-dataset-out

The leave-one-dataset-out protocol evaluates the generalization performance of the trained models by testing them on the test set of an unseen dataset; that is, no images from that dataset are available during training. For each experiment, we hold out the test

⁶The list of discarded images can be found at <https://raysonlaroca.github.io/supp/visapp2022/discarded-images.txt>

set of one dataset as the unseen data, and train every model on all images from the other datasets. As an example, if AOLP’s test set is the current unseen data, the models are trained on all images from Caltech Cars, EnglishLP, UCSD-Stills, ChineseLP, OpenALPR-EU, SSIG-SegPlate, UFPR-ALPR and RodoSol-ALPR, in addition to the images taken from the internet and provided by Laroca et al. (2021b).

We evaluate the models only on the test set images from each unseen dataset, rather than including the training and validation images in the evaluation, so that the results achieved by each model on a given dataset are fully comparable with those achieved by the same model under the traditional-split protocol.

4.4.3 Performance Evaluation

As mentioned in Section 1, in our experiments, the LPs fed to the recognition models were detected using YOLOv4 (Bochkovskiy et al., 2020) – with an input size of 672×416 pixels – rather than cropped directly from the ground truth. This procedure was adopted to better simulate real-world scenarios, as the LPs will not always be detected perfectly, and certain OCR models are not as robust in cases where the region of interest has not been detected so precisely (Gonçalves et al., 2018). We employed the YOLOv4 model for this task because impressive results are consistently being reported in the ALPR context through YOLO-based models (Weihong and Jiaoyang, 2020). Indeed, as detailed in Section 5, YOLOv4 reached an average recall rate above 99.5% in our experiments (we considered as correct the detections with Intersection over Union (IoU) ≥ 0.5 with the ground truth).

For each experiment, we report the number of correctly recognized LPs divided by the number of LPs in the test set. A correctly recognized LP means that all characters on the LP were correctly recognized, as a single incorrectly recognized character can result in the vehicle being incorrectly identified.

Note that the first character in Chinese LPs is a Chinese character that represents the province in which the vehicle is affiliated (Xu et al., 2018; Zhang et al., 2021). Even though Chinese LPs are used in our experiments (see Figure 3d), the evaluated models were not trained/adjusted to recognize Chinese characters; that is, only digits and English letters are considered. This same procedure was adopted in previous works (Li et al., 2019; Selmi et al., 2020; Laroca et al., 2021b) for several reasons, including scope reduction and the fact that it is not trivial for non-Chinese speakers to analyze the different Chinese characters in order to make an accurate error analysis or to choose which data augmentation techniques to explore. Following Li et al. (2019), we denoted all Chinese char-

acters as a single class ‘*’ in our experiments. According to our results, the recognition models learned well the difference between Chinese characters and others – i.e., digits and English letters – and this procedure did not affect the recognition rates obtained.

5 RESULTS AND DISCUSSION

First, we report in Table 4 the recall rates obtained by the YOLOv4 model in the LP detection stage. As can be seen, it reached surprisingly good results in both protocols. More specifically, recall rates above 99.9% were achieved in 14 of the 18 assessments. As in previous works (Laroca et al., 2018; Gonçalves et al., 2018; Silva and Jung, 2020), the detection results are slightly worse for the UFPR-ALPR dataset due to its challenging nature, as (i) it has images where the vehicles are considerably far from the camera; (ii) some of its frames have motion blur because the dataset was recorded in real-world scenarios where both the vehicle and the camera – inside another vehicle – are moving; and (iii) it also contains images of motorcycles, where the backgrounds can be much more complicated due to different body configurations and mixtures with other background scenes (Hsu et al., 2015).

Considering the discussion above, we assert that deep models trained for LP detection on images from multiple datasets can be employed quite reliably on images from unseen datasets (i.e., leave-one-dataset-out protocol). Of course, this may not hold true in extraordinary cases where the test set domain is very different from training ones, but this was not the case in our experimental evaluation carried out on images from nine datasets with different characteristics.

Regarding the recognition stage, the results achieved by all models across all datasets on the traditional-split and leave-one-dataset-out protocols are shown in Table 5 and Table 6, respectively. In Table 6, we included the results obtained by Sighthound (Masood et al., 2017) and OpenALPR (OpenALPR API, 2021), which are two commercial systems frequently used as baselines in the ALPR literature, since in principle they are trained on images from large-scale private datasets and not from the public datasets explored here (i.e., leave-one-dataset-out protocol).

The first observation is that, as expected, the best results – on average for all models – were attained when training and evaluating the models on different subsets from the same datasets (i.e., traditional-split protocol). The only case where this did not occur was precisely in the OpenALPR-EU dataset, where

no images are used for training even under the traditional-split protocol (see Table 3). We kept this division for three main reasons: (i) to better evaluate the recognition models on European LPs; (ii) to maintain consistency with previous works (Masood et al., 2017; Laroca et al., 2021b; Silva and Jung, 2022), which also used all images from that dataset for testing; and (iii) to analyze how the models perform with more training data from other datasets, which in this case corresponds to the leave-one-dataset-out protocol since all images from the other datasets – and not just the training set ones – are used for training. Although some studies have shown that the performance on the test set of a particular dataset often decreases when the training data is augmented with data from other datasets (Torralla and Efros, 2011; Khosla et al., 2012), the recognition rates reached in the OpenALPR-EU dataset were higher with more training data from other datasets. In the same way, CR-NET performed better in the EnglishLP dataset when using all images from the OpenALPR-EU dataset for training (both datasets contain images of European LPs).

The average recognition rate across all datasets decreased from 82.4% under the traditional-split protocol to 74.5% under the leave-one-dataset-out protocol. This drastic performance drop is accentuated by the poor results achieved on the EnglishLP and AOLP datasets under the leave-one-dataset-out protocol. For instance, the average recognition rate of 90.8% obtained in the AOLP dataset under the traditional-split protocol drops to 62.7% under the leave-one-dataset-out protocol. These results caught us by surprise, as both datasets have been considered to be relatively simple due to the fact that recent works have reported recognition rates close to 97% for the EnglishLP and above 99% for the AOLP dataset (Henry et al., 2020; Laroca et al., 2021b; Silva and Jung, 2022; Wang et al., 2022). According to our analysis, most of the recognition errors under the leave-one-dataset-out protocol occurred due to differences in the fonts of the LP characters in the training and test images, as well as because of specific patterns in the LP (e.g., a coat of arms between the LP characters or a straight line under them). To better illustrate, Figure 5 shows three LPs from the AOLP dataset where the TRBA model, which performed best on that dataset, recognized at least one character incorrectly under the leave-one-dataset-out protocol but not under the traditional split. Such analysis highlights the importance of performing cross-dataset experiments in the ALPR context.

The second observation is that, regardless of the evaluation protocol adopted, no recognition model achieved the best results in every single dataset we

Table 4: Recall rates obtained by YOLOv4 in the LP detection stage.

| Approach | Test set | Caltech Cars # 46 | EnglishLP # 102 | UCSD-Stills # 60 | ChineseLP # 161 | AOLP # 687 | OpenALPR-EU # 108 | SSIG-SegPlate # 804 | UFPR-ALPR # 1,800 | RodoSol-ALPR # 8,000 | Average |
|--------------------------------|----------|----------------------|--------------------|---------------------|--------------------|---------------|----------------------|------------------------|----------------------|-------------------------|---------|
| YOLOv4 (traditional-split) | | 100.0% | 100.0% | 100.0% | 100.0% | 100.0% | 100.0% | 99.9% | 99.1% | 100.0% | 99.9% |
| YOLOv4 (leave-one-dataset-out) | | 100.0% | 100.0% | 100.0% | 100.0% | 99.9% | 99.1% | 100.0% | 96.8% | 99.6% | 99.5% |

Table 5: Recognition rates obtained by all models under the traditional-split protocol, which assesses the ability of the models to perform well in seen scenarios. Each model (rows) was trained once on the union of the training set images from all datasets and evaluated on the respective test sets (columns). The best recognition rate achieved in each dataset is shown in bold.

| Approach | Test set | Caltech Cars # 46 | EnglishLP # 102 | UCSD-Stills # 60 | ChineseLP # 161 | AOLP # 687 | OpenALPR-EU # 108 | SSIG-SegPlate # 804 | UFPR-ALPR # 1,800 | RodoSol-ALPR # 8,000 | Average |
|--|----------|----------------------|--------------------|---------------------|--------------------|---------------|----------------------|------------------------|----------------------|-------------------------|--------------|
| CR-NET (Silva and Jung, 2020) | | 95.7% | 92.2% | 100.0% | 96.9% | 97.7% | 97.2% | 97.1% | 78.3% | 55.8% [‡] | 90.1% |
| CRNN (Shi et al., 2017) | | 87.0% | 81.4% | 88.3% | 88.2% | 87.6% | 89.8% | 93.4% | 64.9% | 48.2% | 81.0% |
| Fast-OCR (Laroca et al., 2021a) | | 93.5% | 81.4% | 95.0% | 85.1% | 95.8% | 91.7% | 87.1% | 65.9% | 49.7% [‡] | 82.8% |
| GRCNN (Wang and Hu, 2017) | | 93.5% | 87.3% | 91.7% | 84.5% | 85.9% | 87.0% | 94.3% | 63.3% | 48.4% | 81.7% |
| Holistic-CNN (Špaňhel et al., 2017) | | 89.1% | 68.6% | 88.3% | 90.7% | 86.3% | 78.7% | 94.8% | 70.3% | 49.0% | 79.5% |
| Multi-task (Gonçalves et al., 2019) | | 87.0% | 62.7% | 85.0% | 86.3% | 84.7% | 66.7% | 93.0% | 65.3% | 49.1% | 75.5% |
| R ² AM (Lee and Osindero, 2016) | | 84.8% | 70.6% | 81.7% | 87.0% | 83.1% | 63.9% | 92.0% | 66.9% | 48.6% | 75.4% |
| RARE (Shi et al., 2016) | | 91.3% | 84.3% | 90.0% | 95.7% | 93.4% | 91.7% | 93.7% | 69.0% | 51.6% | 84.5% |
| Rosetta (Borisjuk et al., 2018) | | 87.0% | 75.5% | 81.7% | 90.1% | 83.7% | 81.5% | 94.3% | 63.9% | 48.7% | 84.5% |
| STAR-Net (Liu et al., 2016) | | 95.7% | 93.1% | 96.7% | 96.9% | 96.8% | 95.4% | 96.1% | 70.9% | 51.8% | 88.2% |
| TRBA (Baek et al., 2019) | | 91.3% | 87.3% | 96.7% | 96.9% | 99.0% | 93.5% | 97.3% | 72.9% | 59.6% | 88.3% |
| ViTSTR-Base (Atienza, 2021) | | 84.8% | 80.4% | 90.0% | 99.4% | 95.6% | 84.3% | 96.1% | 73.3% | 49.3% | 83.7% |
| Average | | 90.0% | 80.4% | 90.4% | 91.5% | 90.8% | 85.1% | 94.1% | 68.7% | 50.8% | 82.4% |

[‡]Images from the RodoSol-ALPR dataset were not used for training the CR-NET and Fast-OCR models, as each character’s bounding box needs to be labeled for training them (as detailed in Section 4.1).

Table 6: Recognition rates obtained by all models under the leave-one-dataset-out protocol, which assesses the generalization performance of the models by testing them on the test set of an unseen dataset. For each dataset (columns), we trained the recognition models (rows) on all images from the other datasets. The best recognition rates achieved are shown in bold.

| Approach | Test set | Caltech Cars # 46 | EnglishLP # 102 | UCSD-Stills # 60 | ChineseLP # 161 | AOLP # 687 | OpenALPR-EU # 108 | SSIG-SegPlate # 804 | UFPR-ALPR # 1,800 | RodoSol-ALPR # 8,000 | Average |
|--|----------|----------------------|--------------------|---------------------|--------------------|---------------|----------------------|------------------------|----------------------|-------------------------|--------------|
| CR-NET (Silva and Jung, 2020) | | 93.5% | 96.1% | 96.7% | 88.2% | 76.9% | 96.3% | 94.7% | 61.8% | 45.4% | 83.3% |
| CRNN (Shi et al., 2017) | | 91.3% | 62.7% | 75.0% | 76.4% | 59.4% | 88.0% | 91.3% | 61.7% | 38.8% | 71.6% |
| Fast-OCR (Laroca et al., 2021a) | | 93.5% | 91.2% | 95.0% | 90.1% | 77.0% | 94.4% | 91.2% | 53.2% | 47.8% | 81.5% |
| GRCNN (Wang and Hu, 2017) | | 95.7% | 65.7% | 90.0% | 80.7% | 53.9% | 88.9% | 90.3% | 60.8% | 39.8% | 74.0% |
| Holistic-CNN (Špaňhel et al., 2017) | | 80.4% | 40.2% | 73.3% | 81.4% | 59.7% | 83.3% | 93.4% | 61.8% | 33.4% | 67.4% |
| Multi-task (Gonçalves et al., 2019) | | 82.6% | 34.3% | 66.7% | 77.6% | 50.8% | 79.6% | 89.9% | 57.9% | 44.8% | 64.9% |
| R ² AM (Lee and Osindero, 2016) | | 89.1% | 52.9% | 66.7% | 74.5% | 52.5% | 80.6% | 93.5% | 57.9% | 40.7% | 67.6% |
| RARE (Shi et al., 2016) | | 84.8% | 50.0% | 85.0% | 88.8% | 62.9% | 91.7% | 93.5% | 71.3% | 40.1% | 74.2% |
| Rosetta (Borisjuk et al., 2018) | | 89.1% | 63.7% | 68.3% | 83.2% | 51.1% | 81.5% | 94.4% | 61.8% | 42.5% | 70.6% |
| STAR-Net (Liu et al., 2016) | | 89.1% | 80.4% | 91.7% | 95.0% | 79.3% | 93.5% | 94.0% | 69.1% | 43.6% | 81.8% |
| TRBA (Baek et al., 2019) | | 95.7% | 66.7% | 93.3% | 95.0% | 70.0% | 92.6% | 97.3% | 96.9% | 42.6% | 80.7% |
| ViTSTR-Base (Atienza, 2021) | | 89.1% | 58.8% | 90.0% | 95.0% | 59.2% | 89.8% | 97.9% | 69.6% | 41.7% | 76.8% |
| Average | | 89.5% | 63.6% | 82.6% | 85.5% | 62.7% | 88.3% | 93.4% | 63.3% | 41.8% | 74.5% |
| Average (traditional-split protocol) | | 90.0% | 80.4% | 90.4% | 91.5% | 90.8% | 85.1% [†] | 94.1% | 68.7% | 50.8% | 82.4% |
| Sighthound (Masood et al., 2017) | | 87.0% | 94.1% | 90.0% | 84.5% | 79.6% | 94.4% | 79.2% | 52.6% | 51.0% | 79.2% |
| OpenALPR (OpenALPR API, 2021)* | | 95.7% | 99.0% | 96.7% | 93.8% | 81.1% | 99.1% | 91.4% | 87.8% | 70.0% | 90.5% |

[†]Under the traditional-split protocol, no images from the OpenALPR-EU dataset were used for training. This is the protocol commonly adopted in the literature (Laroca et al., 2021b; Silva and Jung, 2022).

*OpenALPR contains specialized solutions for LPs from different regions and the user must enter the correct region before using its API. Hence, it was expected to achieve better results than the other methods.

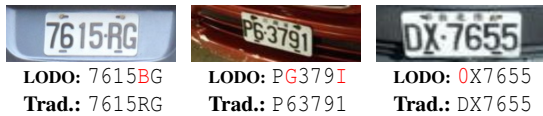


Figure 5: The predictions obtained by TRBA on three images of the AOLP dataset. In general, the errors (outlined in red) under the leave-one-dataset-out (LODO) protocol did not occur in challenging cases (e.g., blurry or tilted images); therefore, they were probably caused by *differences in the training and test images*. Trad.: traditional-split protocol.

performed experiments on. For instance, although the CR-NET model obtained the best average recognition rates, corroborating the state-of-the-art results reported recently in (Laroca et al., 2021b; Silva and Jung, 2022), it did not reach the best results in the ChineseLP, AOLP, SSIG-SegPlate and RodoSol-ALPR datasets in either protocol. These results emphasize the importance of carrying out experiments on multiple datasets, with LPs from

different countries/regions, especially under the leave-one-dataset-out protocol because six different models obtained the best result in at least one dataset.

The third observation is that the RodoSol-ALPR dataset proved very challenging since all the recognition models trained by us, as well as both commercial systems, failed to reach recognition rates above 70% on its test set images. The main reason for such disappointing results is the large number of motorcycle images, which are very challenging in nature (as discussed in Section 2). For example, OpenALPR correctly recognized 3,772 of the 4,000 cars in the test set (94.3%) and only 1,827 of the 4,000 motorcycles in the test set (45.7%). These results accentuate the importance of the proposed dataset for the accurate evaluation of ALPR systems, as it avoids bias in the assessments by having the same number of “easy” (cars with single-row LPs) and “difficult” (motorcycles with two-row LPs) samples.

We also did not rule out challenging images when selecting the images for the creation of the dataset. Figure 6 shows some of these images along with the predictions returned by TRBA (traditional-split) and OpenALPR, which were the model trained by us and the commercial system that performed better on this dataset. The results are in line with what was recently stated by Zhang et al. (2021): that recognizing LPs in complex environments is still far from satisfactory.



Figure 6: Some LP images from RodoSol-ALPR along with the predictions returned by TRBA and OpenALPR. Note that one character may become very similar to another due to factors such as blur, low/high exposure, rotations and occlusions. For correctness, we checked if the ground truth (GT) matched the vehicle make and model on the National Traffic Department of Brazil (DENATRAN) database.

Lastly, it is important to highlight the number of experiments we carried out for this work. We trained each of the 12 chosen OCR models 10 times: once following the split protocols traditionally adopted in the literature (see Table 5) and nine for the leave-one-dataset-out evaluation (see Table 6). We remark that a single training process of some models (e.g., TRBA and ViTSTR-Base) took several days to complete on an NVIDIA Quadro RTX 8000 GPU. In fact, we believe that this large number of necessary experiments is precisely what caused a leave-one-dataset-out evaluation to have not yet been performed in the literature.

6 CONCLUSIONS

As the performance of traditional-split LP recognition is rapidly improving, researchers should pay more attention to cross-dataset LP recognition since it better simulates real-world ALPR applications, where new cameras are regularly being installed in new locations without existing systems being retrained every time.

As a first step towards that direction, in this work we evaluated 12 OCR models for LP recognition on 9 public datasets with a great variety in several aspects (e.g., acquisition settings, image resolution, and LP layouts). We adopted a traditional-split *versus* leave-one-dataset-out experimental setup to empirically as-

sess the cross-dataset generalization of the chosen models. It is noteworthy that we are not aware of any work in the ALPR context where so many methods were implemented and compared or where so many datasets were explored in the experiments.

As expected, the experimental results showed significant drops in performance for most datasets when training and testing the recognition models in a leave-one-dataset-out fashion. The fact that very low recognition rates (around 63%) were reported in the EnglishLP and AOLP datasets underscored the importance of carrying out cross-dataset experiments, as very high recognition rates (above 95% and 99%, respectively) are frequently achieved on these datasets under the traditional-split protocol.

The importance of exploring various datasets in the evaluation was also demonstrated, as no model performed better than the others in all experiments. It was quite unexpected for us that six different models reached the best result in at least one dataset under the leave-one-dataset-out protocol. In this sense, we draw attention to the fact that most works in the literature used three or fewer datasets in the experiments, although this has been gradually changing in recent years (Selmi et al., 2020; Laroca et al., 2021b).

We also introduced a publicly available dataset for ALPR that, to the best of our knowledge, is the first to contain images of vehicles with Mercosur LPs. We expect it will assist in developing new approaches for this LP layout and the fair comparison between methods proposed in different works. Additionally, the proposed dataset includes 10,000 motorcycle images, being by far the largest in this regard. RodoSol-ALPR has proved challenging in our experiments, as both the models trained by us and two commercial systems reached recognition rates below 70% on its test set.

As future work, we plan to gather images from the internet to build a novel dataset for end-to-end ALPR with images acquired in various countries/regions, by many different cameras, both static or mobile, with a well-defined evaluation protocol for both intra- and cross-dataset LP detection and LP recognition. In addition, we intend to leverage the potential of Generative Adversarial Networks (GANs) to generate hundreds of thousands of synthetic LP images with different transformations and balanced character classes in order to improve the generalization ability of deep models. Finally, we would like to carry out more experiments to quantify the influence of each dataset, especially RodoSol-ALPR, on the performance of the models under the leave-one-dataset-out protocol.

ACKNOWLEDGMENTS

This work was supported in part by the Coordination for the Improvement of Higher Education Personnel (CAPES) (Social Demand Program), and in part by the National Council for Scientific and Technological Development (CNPq) (Grant 308879/2020-1). The Quadro RTX 8000 GPU used for this research was donated by the NVIDIA Corporation. We also thank the *Rodovia do Sol* (RodoSol) concessionaire, particularly the information technology (IT) manager Marciano Calvi Ferri, for providing the images for the creation of the RodoSol-ALPR dataset.

REFERENCES

- Ashraf, A., Khan, S. S., Bhagwat, N., and Taati, B. (2018). Learning to unlearn: Building immunity to dataset bias in medical imaging studies. In *Machine Learning for Health Workshop at NeurIPS 2018*, pages 1–5.
- Atienza, R. (2021). Vision transformer for fast and efficient scene text recognition. In *International Conference on Document Analysis and Recognition*, pages 319–334.
- Baek, J. et al. (2019). What is wrong with scene text recognition model comparisons? dataset and model analysis. In *IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 4714–4722.
- Beratoglu, M. S. and Toreyin, B. U. (2021). Vehicle license plate detector in compressed domain. *IEEE Access*, 9:95087–95096.
- Bochkovskiy, A., Wang, C.-Y., and Liao, H.-Y. M. (2020). YOLOv4: Optimal speed and accuracy of object detection. *arXiv preprint*, arXiv:2004.10934:1–14.
- Borisyuk, F., Gordo, A., and Sivakumar, V. (2018). Rosetta: Large scale system for text detection and recognition in images. In *ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, page 71–79.
- Buslaev, A. et al. (2020). Albuementations: Fast and flexible image augmentations. *Information*, 11(2):125.
- Dlagnekov, L. and Belongie, S. (2005). UCSD dataset. https://www.bongielab.org/car_data.html.
- Estevam, V., Laroca, R., Pedrini, H., and Menotti, D. (2021). Tell me what you see: A zero-shot action recognition method based on natural language descriptions. *arXiv preprint*, arXiv:2112.09976:1–15.
- Forbes (2021). Poor Auto Financials Likely as Sales Sag, but Forecasts Point to Strong Turnaround. <https://www.forbes.com/sites/neilwinton/2021/10/10/poor-auto-financials-likely-as-sales-sag-but-forecasts-point-to-strong-turnaround/>.
- Gonçalves, G. R., da Silva, S. P. G., Menotti, D., and Schwartz, W. R. (2016). Benchmark for license plate character segmentation. *Journal of Electronic Imaging*, 25(5):053034.
- Gonçalves, G. R., Diniz, M. A., Laroca, R., Menotti, D., and Schwartz, W. R. (2019). Multi-task learning for low-resolution license plate recognition. In *Iberoamerican Congress on Pattern Recognition*, pages 251–261.
- Gonçalves, G. R. et al. (2018). Real-time automatic license plate recognition through deep multi-task networks. In *Conference on Graphics, Patterns and Images (SIB-GRAPI)*, pages 110–117.
- Hasnat, A. and Nakib, A. (2021). Robust license plate signatures matching based on multi-task learning approach. *Neurocomputing*, 440:58–71.
- Henry, C., Ahn, S. Y., and Lee, S.-W. (2020). Multinational license plate recognition using generalized character sequence detection. *IEEE Access*, 8:35185–35199.
- Hsu, G., Zeng, S., Chiu, C., and Chung, S. (2015). A comparison study on motorcycle license plate detection. In *IEEE International Conference on Multimedia Expo Workshops (ICMEW)*, pages 1–6.
- Hsu, G. S., Chen, J. C., and Chung, Y. Z. (2013). Application-oriented license plate recognition. *IEEE Trans. on Vehicular Technology*, 62(2):552–561.
- IHS Markit (2021). 2022 global light vehicle production outlook intact. <https://ihsmarkit.com/research-analysis/2022-global-light-vehicle-production-outlook.html>.
- Izidio, D. M. F. et al. (2020). An embedded automatic license plate recognition system using deep learning. *Design Automation for Embedded Systems*, 24:23–43.
- Kessentini, Y., Besbes, M. D., Ammar, S., and Chabbouh, A. (2019). A two-stage deep neural network for multi-norm license plate detection and recognition. *Expert Systems with Applications*, 136:159–170.
- Khosla, A., Zhou, T., Malisiewicz, T., Efros, A. A., and Torralba, A. (2012). Undoing the damage of dataset bias. In *European Conference on Computer Vision (ECCV)*, pages 158–171.
- Laroca, R., Araujo, A. B., Zanlorensi, L. A., de Almeida, E. C., and Menotti, D. (2021a). Towards image-based automatic meter reading in unconstrained scenarios: A robust and efficient approach. *IEEE Access*, 9:67569–67584.
- Laroca, R., Barroso, V., Diniz, M. A., Gonçalves, G. R., Schwartz, W. R., and Menotti, D. (2019). Convolutional neural networks for automatic meter reading. *Journal of Electronic Imaging*, 28(1):013023.
- Laroca, R. et al. (2018). A robust real-time automatic license plate recognition based on the YOLO detector. In *International Joint Conference on Neural Networks (IJCNN)*, pages 1–10.
- Laroca, R., Zanlorensi, L. A., Gonçalves, G. R., Todt, E., Schwartz, W. R., and Menotti, D. (2021b). An efficient and layout-independent automatic license plate recognition system based on the YOLO detector. *IET Intelligent Transport Systems*, 15(4):483–503.
- Lee, C. and Osindero, S. (2016). Recursive recurrent nets with attention modeling for OCR in the wild. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 2231–2239.
- Li, H., Wang, P., and Shen, C. (2019). Toward end-to-end car license plate detection and recognition with deep neural networks. *IEEE Transactions on Intelligent Transportation Systems*, 20(3):1126–1136.

- Liang, J. et al. (2021). EGSA-Net: edge-guided sparse attention network for improving license plate detection in the wild. *Applied Intelligence*, 52(4):4458–4472.
- Liu, W., Chen, C., Kwan-Yee K. Wong, Z. S., and Han, J. (2016). STAR-Net: A spatial attention residue network for scene text recognition. In *British Machine Vision Conference (BMVC)*, pages 1–13.
- Lubna, Mufti, N., and Shah, S. A. A. (2021). Automatic number plate Recognition: A detailed survey of relevant algorithms. *Sensors*, 21(9):3028.
- Masood, S. Z. et al. (2017). License plate detection and recognition using deeply learned convolutional neural networks. *arXiv preprint*, arXiv:1703.07330.
- Meng, S., Zhang, Z., and Wan, Y. (2020). Accelerating automatic license plate detection in the wild. In *IEEE Joint International Information Technology and Artificial Intelligence Conference*, pages 742–746.
- Oliveira, I. O. et al. (2021). Vehicle-Rear: A new dataset to explore feature fusion for vehicle identification using convolutional neural networks. *IEEE Access*, 9:101065–101077.
- OpenALPR API (2021). <http://www.openalpr.com/>.
- OpenALPR Inc. (2016). OpenALPR-EU dataset. <https://github.com/openalpr/benchmarks/tree/master/endoend/eu>.
- Panahi, R. and Gholampour, I. (2017). Accurate detection and recognition of dirty vehicle plate numbers for high-speed applications. *IEEE Transactions on Intelligent Transportation Systems*, 18(4):767–779.
- Presidência da República (1997). LEI N° 9.503, DE 23 DE SETEMBRO DE 1997 - Código de Trânsito Brasileiro. http://www.planalto.gov.br/ccivil_03/leis/19503compilado.htm.
- Redmon, J., Divvala, S., Girshick, R., and Farhadi, A. (2016). You only look once: Unified, real-time object detection. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 779–788.
- RodoSol (2022). *Concessionária Rodovia do Sol S/A*. <https://www.rodosol.com.br/blog/conheca-ardosol-2>. Accessed: 2022-02-06.
- Selmi, Z., Halima, M. B., Pal, U., and Alimi, M. A. (2020). DELP-DAR system for license plate detection and recognition. *Pattern Recog. Letters*, 129:213–223.
- Shi, B., Bai, X., and Yao, C. (2017). An end-to-end trainable neural network for image-based sequence recognition and its application to scene text recognition. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 39(11):2298–2304.
- Shi, B., Wang, X., Lyu, P., Yao, C., and Bai, X. (2016). Robust scene text recognition with automatic rectification. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 4168–4176.
- Silva, S. M. and Jung, C. R. (2020). Real-time license plate detection and recognition using deep convolutional neural networks. *Journal of Visual Communication and Image Representation*, page 102773.
- Silva, S. M. and Jung, C. R. (2022). A flexible approach for automatic license plate recognition in unconstrained scenarios. *IEEE Transactions on Intelligent Transportation Systems*, 23(6):5693–5703.
- Silvano, G. et al. (2021). Synthetic image generation for training deep learning-based automated license plate recognition systems on the Brazilian Mercosur standard. *Design Automation for Embedded Systems*, 25(2):113–133.
- Špaňhel, J. et al. (2017). Holistic recognition of low quality license plates by CNN using track annotated data. In *IEEE International Conference on Advanced Video and Signal Based Surveillance (AVSS)*, pages 1–6.
- Srebrić, V. (2003). EnglishLP database. https://www.zemris.fer.hr/projects/LicensePlates/english/baza_slika.zip.
- Tommasi, T., Patricia, N., Caputo, B., and Tuytelaars, T. (2017). A deeper look at dataset bias. In *Domain Adaptation in Computer Vision Applications*, pages 37–55. Springer.
- Torralba, A. and Efros, A. A. (2011). Unbiased look at dataset bias. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 1521–1528.
- Wang, J. and Hu, X. (2017). Gated recurrent convolution neural network for OCR. In *Annual Conference on Neural Information Processing Systems (NeurIPS)*.
- Wang, Y. et al. (2022). Rethinking and designing a high-performing automatic license plate recognition approach. *IEEE Transactions on Intelligent Transportation Systems*, 23(7):8868–8880.
- Weber, M. (1999). Caltech Cars dataset. <https://data.caltech.edu/records/20084>.
- Weihong, W. and Jiaoyang, T. (2020). Research on license plate recognition algorithms based on deep learning in complex environment. *IEEE Access*, 8:91661–91675.
- Xiang, H., Zhao, Y., Yuan, Y., Zhang, G., and Hu, X. (2019). Lightweight fully convolutional network for license plate detection. *Optik*, 178:1185–1194.
- Xie, L., Ahmad, T., Jin, L., Liu, Y., and Zhang, S. (2018). A new CNN-based method for multi-directional car license plate detection. *IEEE Transactions on Intelligent Transportation Systems*, 19(2):507–517.
- Xu, Z. et al. (2018). Towards end-to-end license plate detection and recognition: A large dataset and baseline. In *European Conf. on Computer Vision*, pages 261–277.
- Zhang, J., Li, W., Ogunbona, P., and Xu, D. (2019). Recent advances in transfer learning for cross-dataset visual recognition: A problem-oriented perspective. *ACM Computing Surveys*, 52(1):1–38.
- Zhang, L., Wang, P., Li, H., Li, Z., Shen, C., and Zhang, Y. (2021). A robust attentional framework for license plate recognition in the wild. *IEEE Trans. on Intelligent Transportation Systems*, 22(11):6967–6976.
- Zhou, W., Li, H., Lu, Y., and Tian, Q. (2012). Principal visual word discovery for automatic license plate detection. *IEEE Transactions on Image Processing*, 21(9):4269–4279.
- Zou, Y., Zhang, Y., Yan, J., Jiang, X., Huang, T., Fan, H., and Cui, Z. (2020). A robust license plate recognition model based on Bi-LSTM. *IEEE Access*, 8:211630.