# A TWO-STREAM SIAMESE NEURAL NETWORK FOR VEHICLE RE-IDENTIFICATION BY USING NON-OVERLAPPING CAMERAS

*Icaro O. de Oliveira, Keiko V. O. Fonseca and Rodrigo Minetto*

Federal University of Technology - Paraná (UTFPR), Curitiba, Brazil
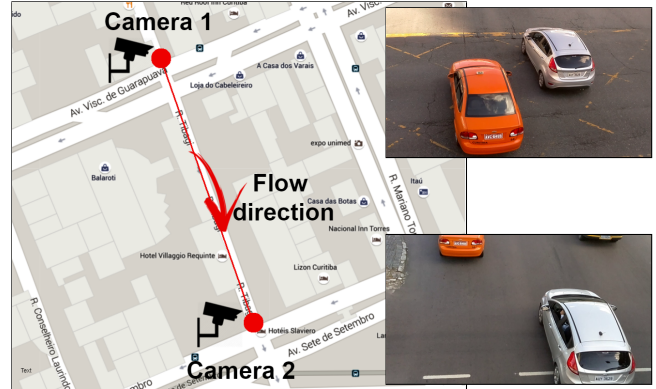
## ABSTRACT

We describe in this paper a Two-Stream Siamese Neural Network for vehicle re-identification. The proposed network is fed simultaneously with small coarse patches of the vehicle shape's, with $96 \times 96$ pixels, in one stream, and fine features extracted from license plate patches, easily readable by humans, with $96 \times 48$ pixels, in the other one. Then, we combined the strengths of both streams by merging the Siamese distance descriptors with a sequence of fully connected layers, as an attempt to tackle a major problem in the field, false alarms caused by a huge number of car design and models with nearly the same appearance or by similar license plate strings. In our experiments, with 2 hours of videos containing 2982 vehicles, extracted from two low-cost cameras in the same roadway, 546 ft away, we achieved a $F$-measure and accuracy of 92.6% and 98.7%, respectively. We show that our network, available at https://github.com/icarofua/siamese-two-stream, outperforms other One-Stream architectures, even if they use higher resolution image features.

***Index Terms***— Vehicle Re-identification; Siamese Neural Networks; Vehicle Matching; Travel Time Estimation.

## 1. INTRODUCTION

This paper address the problem of matching moving vehicles that appear in two videos taken by multiple cameras with non-overlapping fields of view. See Fig. 1. This is a common sub-problem for several applications in intelligent transportation systems, such as enforcement of road speed limits, criminal investigations, monitoring of commercial transportation vehicles, and traffic management.

Some of these applications traditionally use physical sensors placed over, near, or under the road, such as pressure-sensitive cables and inductive loop detectors [1, 2]. However, these detectors present limitations such as the same vehicles entering or leaving the road between the two cameras. Other applications use optical character recognition (OCR)

Icaro O. de Oliveira, Keiko V. O. Fonseca and Rodrigo Minetto are with the Graduate Program in Electrical and Computer Engineering (CPGEI) and PPGCA. E-mail: icarofua@gmail.com

**Fig. 1**. System setup: a traffic engineering company placed in two different semaphores a pair of low-cost full-HD cameras properly calibrated and time synchronized. In general, not every vehicle seen in one video appears in the other video.

algorithms [3] to translate the license plate image regions into character codes, such as ASCII. However, this translation is not straightforward when two or more lanes are recorded at the same time, producing small license plate regions that are very hard to read. Recognition of vehicles by shape and color is not sufficiently reliable either, since vehicles of the same brand and model often look exactly the same [4].

For such reasons, in our solution we have opted to identify vehicles across non-overlapping cameras by using an hybrid strategy, that is, we developed a Two-Stream Siamese Neural Network that is fed, simultaneously, with two of the most distinctive and persistent features available, the vehicle's shape and the registration license plate. Then, for fusion of the Two-Streams we concatenate the distance descriptors extracted from each single Siamese network and add fully connected layers for classification. We also show that the combination of small image patches produces a fast network that outperforms other complex architectures, even if they use higher resolution image patches. The rest of this paper is organized as follows. In Sec. 2, we discuss the related work. In Sec. 3, we describe the Two-Stream Siamese Network. Experiments are reported in Sec. 4. Finally, in Sec. 5 we state the conclusions.

## 2. RELATED WORK

Vehicle re-identification is an active field of research with many algorithms and extensive bibliography [1, 2, 5, 6, 7, 8, 9]. The survey of Tian *et al.* [10] listed this problem as an open challenge for intelligent transportation systems. Traditionally, algorithms for this task were based on the comparison of electromagnetic signatures. However, as observed by Ndoye *et al.* [2], such signature-matching algorithms are exceedingly complex and depends on extensive calibrations or complicated data models.
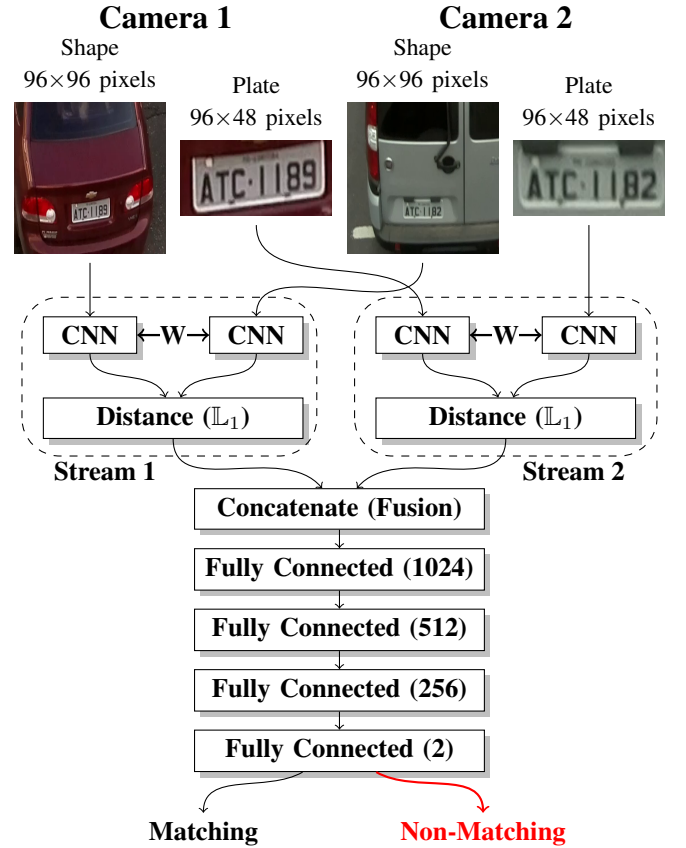
Video-based algorithms have been proven to be powerful for vehicle re-identification [2, 7, 8, 9, 11]. Such algorithms need to address *fine-grained vehicle recognition* issues [12], that is, to distinguish between subordinate categories with similar visual appearance, caused by a huge number of car design and models with similar appearance. As an attempt to solve these issues many authors proposed to use handed-crafted image descriptors such as SIFT [13]. Recently, inspired by the tremendous progress of the Siamese Neural Networks Tang *et al.* [9], in 2017, proposed for vehicle re-identification in traffic surveillance environment to fuse deep and hand-crafted features by using a Siamese Triplet Network [14]. In 2018, Yan *et al.* [7] proposed a novel deep learning metric, a Triplet Loss Function, that takes into account the inter-class similarity and intra-class variance in vehicle models considering only the vehicle's shape. Also in 2018, Liu *et al.* [8] proposed a coarse-to-fine vehicle re-identification algorithm that initially filters out the potential matchings by using hand-crafted and deep features based on shape and color and, then they used the license plates in a Siamese Network and a Spatiotemporal re-ranking to refine the search.

The idea of a two stream convolutional neural networks (CNN) is not new. Ye *et al.* [15] proposed an architecture that uses static video frames as input in one stream and optical flow features in the other stream for video classification. Chung *et al.* [16] also proposed a two-stream architecture composed by two Siamese CNN fed with spatial and temporal information extracted from RGB frames and optical flow vectors for person re-identification. Zagoruyko *et al.* [17] described distinct Siamese architectures to compare learning image patches. In special, the Central-Surround Two-Stream architecture is similar to the one proposed here.

Finally, some authors [1] use self-adaptive time-window constraints to define upper and lower bounds in order to predict the search space and narrow-down the potential matches. That is, they predict a time-window size based on the camera's distance and the traffic conditions, e.g. free flow or congested. However, we are not trying to solve the travel time estimation problem here, thus, we considered the maximum number of true or false matchings available in order to evaluate the robustness of the architectures.
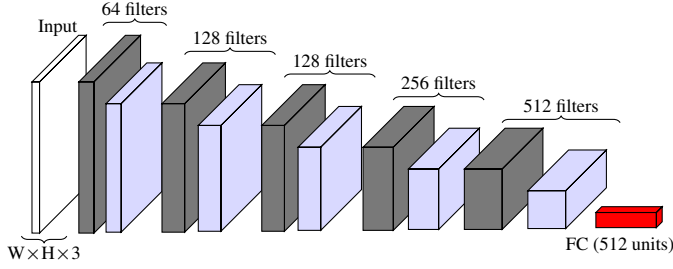
## 3. TWO-STREAM SIAMESE NETWORK

The inference flowchart of the proposed Two-Stream Siamese Network is shown in Fig. 2. The left stream processes the vehicle's shapes while the right stream the license plates. The network weights $W$ are shared only within each stream. We merged the distance vectors of each Siamese — whose similarity is measured by a Mahalanobis distance — and combined the strengths of both features by using a sequence of fully connected layers with dropout regularization (20%) in order to avoid over-fitting. Then, we used a softmax activation function to classify matching pairs from non-matching pairs.



**Fig. 2**. Inference flowchart of the proposed Two-Stream Siamese for Vehicle Matching.

We extracted the vehicle rear end and the vehicle license plate by using the real-time motion detector and algorithms described by Luvizon *et al.* [18, 19]. The CNN used in our Siamese is shown in Fig. 3. Basically, it is a simplified VGG [20] based network, with a reduced number of layers so as to save computational effort. Each CNN provided a vector with 512 features. Each Distance ($\mathbb{L}_1$) provided a vector with 512 distances. Finally, Concatenate (Fusion) provided a vector with 1024 distances.

**Fig. 3**. Small-VGG: a VGG-based convolutional neural network used in the Two-Stream Siamese Network. The dark gray boxes denote convolutions by using filter kernel size of $3 \times 3$; the light blue boxes denote $2 \times 2$ max-pooling layers; and, the red box depicts a fully connected layer.



**Fig. 4**. Pair-wise data augmentation only for positive (true) matchings: from $N$ images of two vehicle sequences we extract $N^2$ distinct pairs. The same procedure is applied for license plate and vehicle shapes.
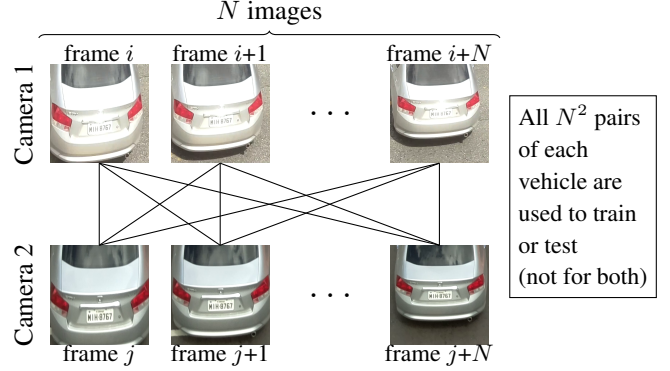
## 4. EXPERIMENTS

For our tests, we used 10 videos — 5 from Camera 1 and 5 from Camera 2 (20 minutes of duration each one) — recorded with frame resolution of $1920 \times 1080$ pixels, at 30.15 frames per second. They are summarized in Table 1.

**Table 1**. Dataset information: number of vehicles and number of vehicles with a visible license plate in Camera 1 and 2; number of vehicles matchings between Camera 1 and 2.

| Set | Camera 1 | | Camera 2 | | No.Match. |
|---|---|---|---|---|---|
| | #Vehicles | #Plates | #Vehicles | #Plates | |
| 01 | 389 | 343 | 280 | 245 | 199 |
| 02 | 350 | 310 | 244 | 227 | 174 |
| 03 | 340 | 301 | 274 | 248 | 197 |
| 04 | 280 | 251 | 233 | 196 | 140 |
| 05 | 345 | 295 | 247 | 194 | 159 |
| Total | 1704 | 1500 | 1278 | 1110 | 869 |

There are multiple distinct occurrences of the same vehicle as it moves across the video. Therefore, instead of only 869 matchings as shown in Table 1, we can generate thousands of true matchings by doing the Cartesian product between a sequence of images of the same vehicle that appears in Camera 1 and 2. This data augmentation is usually necessary for CNN training. Therefore, we used the MOSSE tracker [21] to extract the $N$ first occurrences of each license plate (see Fig 4). Note that negative pairs are easier to generate, since we can use any combination of distinct vehicles from Camera 1 and 2.

We also adjusted another parameter $\lambda$ that was meant to multiply the number of false negatives pairs (non-matchings) in the testing set to simulate the network in a real environment assuming it may have many more tests of non-matchings pairs than the opposite. In Table 2 we show some parameter settings for our experiments. Note however, that we keep

the same proportion of positive and negative pairs during the training in order to avoid class imbalance.

**Table 2**. Parameter settings used in our experiments.

| Settings | Training | | Testing | |
|---|---|---|---|---|
| | #positives | #negatives | #positives | #negatives |
| $N = 3, \lambda = 5$ | 3867 | 3867 | 3903 | 19515 |
| $N = 10, \lambda = 10$ | 42130 | 42130 | 42707 | 427070 |

The quantitative criteria we used to evaluate the architectures performance are the precision $P$, recall $R$, accuracy $A$ and $F$-measure. As shown in Table 3, the Two-Stream Siamese outperforms two distinct One-Stream Siamese Networks: the first one, Siamese-Car, is fed only with the shape of vehicles ($96 \times 96$ pixels); the second, Siamese-Plate, only use patches of license plates ($96 \times 48$ pixels). Note that even when we increased the number of false matchings in the negative testing set, $\lambda = 10$, the $F$-measure of the Two-Stream Siamese was similar in both scenarios. The accuracy $A$ is usually much higher since the number of negative pairs is much larger. Some inference results are shown in Fig. 6.

We also tried different CNN in our Two-Stream Siamese, their performance are reported in Table 4. Furthermore, as can be seen in Fig. 5, we also evaluated the performance of the proposed Two-Stream Siamese against two One-Stream Siamese versions fed with larger image patches ($224 \times 224$ pixels). Note that we achieved a higher $F$-measure by using two small image patches than a single patch containing both features. Another advantage is the Two-Stream Siamese training time: 1938 seconds per epoch ($N = 10$ and $\lambda = 10$) against 3441 seconds per epoch of the Siamese-Car by using the same Small-VGG and 4937 seconds with ResNet. The experiments were carried out on a Intel i7 with 32GB DRAM and a Nvidia Titan Xp GPU.

3

**Table 3**. Matching performance of the proposed Two-Stream Siamese (Small-VGG) against two One-Stream Siamese (Car and Plate with Small-VGG) by using different settings to generate image pairs.

| | $N = 3, \lambda = 5$ | | | |
|---|---|---|---|---|
| Algorithm | $P$ | $R$ | $F$ | $A$ |
| Siamese-Car (Stream 1) | 85.8% | 93.1% | 89.3% | 96.3% |
| Siamese-Plate (Stream 2) | 75.9% | 81.8% | 78.8% | 92.6% |
| Siamese (Two-Stream) | 92.7% | 93.0% | 92.9% | 97.6% |
| | $N = 10, \lambda = 10$ | | | |
| Algorithm | $P$ | $R$ | $F$ | $A$ |
| Siamese-Car (Stream 1) | 92.4% | 83.5% | 87.8% | 97.9% |
| Siamese-Plate (Stream 2) | 86.8% | 59.5% | 70.6% | 95.5% |
| Siamese (Two-Stream) | 94.7% | 90.6% | 92.6% | 98.7% |

**Table 4**. Matching performance of the proposed Two-Stream Siamese with different CNN architectures.

| | $N = 10, \lambda = 10$ | | | |
|---|---|---|---|---|
| Siamese (Two-Stream) | $P$ | $R$ | $F$ | $A$ |
| CNN = Lenet5 | 89.6% | 85.2% | 87.3% | 97.8% |
| CNN = Matchnet [22] | 94.5% | 87.1% | 90.7% | 98.4% |
| CNN = MC-CNN [23] | 89.0% | 90.1% | 89.6% | 98.1% |
| CNN = GoogleNet | 88.8% | 81.8% | 85.1% | 97.4% |
| CNN = AlexNet | 91.3% | 86.5% | 88.8% | 98.0% |
| CNN = Small-VGG | 94.7% | 90.6% | 92.6% | 98.7% |

Vehicle (96×96 pixels)   Vehicle (patches 224×224 pixels)



**VS**

Plate (96×48 pixels)

Siamese Two-Stream (**Small-VGG**)
$F$ = 92.6% and $A$ 98.7%

Siamese-Car (**Small-VGG**):
$F$ = 88.1% and $A$ = 97.9%
Siamese-Car (**Resnet50**):
$F$ = 81.2% and $A$ = 97.1%

**Fig. 5**. Siamese Two-Stream versus Siamese-Car.

## 5. CONCLUSIONS

We proposed in this paper a fast Two-Stream Siamese that combines the discriminatory power of two distinctive and persistent features, the vehicle's shape and the registration



Siamese-Car (Stream 1): **non-matching** ✗
Siamese-Plate (Stream 2): **non-matching** ✗
Siamese (Two-Stream): **non-matching** ✗

Siamese-Car (Stream 1): **matching** ✗
Siamese-Plate (Stream 2): **non-matching** ✓
Siamese (Two-Stream): **non-matching** ✓

Siamese-Car (Stream 1): **non-matching** ✓
Siamese-Plate (Stream 2): **matching** ✗
Siamese (Two-Stream): **non-matching** ✓

Siamese-Car (Stream 1): **matching** ✓
Siamese-Plate (Stream 2): **matching** ✓
Siamese (Two-Stream): **matching** ✓

**Fig. 6**. Inference results (testing set): from top-to-bottom an example where the three architectures failed (severe lighting conditions); Siamese-Car failed (similar vehicle shape); Siamese-Plate failed (similar license plate); and, at bottom, the three architectures found a correct matching.

plate, to address the problem of vehicle re-identification by using non-overlapping cameras. Tests indicate that our network is more robust than other One-Stream Siamese architectures which are fed with the same features or larger images. We also evaluated simple and complex CNNs, used by the Siamese Network, to find a trade-off between efficiency and performance.

# 6. REFERENCES

[1] W. H. Lin and D. Tong, "Vehicle re-identification with dynamic time windows for vehicle passage time estimation," *IEEE Trans. on Intelligent Transportation Systems (ITS)*, vol. 12, no. 4, pp. 1057–1063, 2011.

[2] M. Ndoye, V. F. Totten, J. V. Krogmeier, and D. M. Bullock, "Sensing and signal processing for vehicle re-identification and travel time estimation," *IEEE Trans. ITS*, vol. 12, no. 1, pp. 119–131, March 2011.

[3] Jiwei Li, Minh-Thang Luong, and Dan Jurafsky, "A hierarchical neural autoencoder for paragraphs and documents," *CoRR*, vol. abs/1506.01057, 2015.

[4] Kai She, George Bebis, Haisong Gu, and Ronald Miller, "Vehicle tracking using on-line fusion of color and shape features," in *IEEE Conf. on Intelligent Transportation Systems*, 2004, pp. 731–736.

[5] Xian Zhong, Meng Feng, Wenxin Huang, Zheng Wang, and Shinichi Satoh, "Poses guide spatiotemporal model for vehicle re-identification," in *International Conference on Multimedia Modeling*. Springer, 2019, pp. 426–439.

[6] Y. Li, Y. Li, H. Yan, and J. Liu, "Deep joint discriminative learning for vehicle re-identification and retrieval," in *IEEE International Conference on Image Processing (ICIP)*, Sep. 2017, pp. 395–399.

[7] Y. Bai, Y. Lou, F. Gao, S. Wang, Y. Wu, and L. Duan, "Group sensitive triplet embedding for vehicle re-identification," *IEEE Trans. on Multimedia*, pp. 1–1, 2018.

[8] X. Liu, W. Liu, T. Mei, and H. Ma, "Provid: Progressive and multimodal vehicle reidentification for large-scale urban surveillance," *IEEE Transactions on Multimedia*, vol. 20, no. 3, pp. 645–658, March 2018.

[9] Yi Tang, Di Wu, Zhi Jin, Wenbin Zou, and Xia Li, "Multi-modal metric learning for vehicle re-identification in traffic surveillance environment," in *IEEE Int. Conference on Image Processing (ICIP)*, 2017, pp. 2254–2258.

[10] B. Tian, B. T. Morris, M. Tang, Y. Liu, Y. Yao, C. Gou, D. Shen, and S. Tang, "Hierarchical and networked vehicle surveillance in ITS: A survey," *IEEE Trans. on Intelligent Transportation Systems (ITS)*, vol. 16, no. 2, pp. 557–580, April 2015.

[11] N. Jiang, Y. Xu, Z. Zhou, and W. Wu, "Multi-attribute driven vehicle re-identification with spatial-temporal re-ranking," in *IEEE International Conference on Image Processing (ICIP)*, 2018, pp. 858–862.

[12] L. Liao, R. Hu, J. Xiao, Q. Wang, J. Xiao, and J. Chen, "Exploiting effects of parts in fine-grained categorization of vehicles," in *IEEE International Conference on Image Processing (ICIP)*, Sep. 2015, pp. 745–749.

[13] Changyou Zhang, Xiaoya Wang, Jun Feng, Yu Cheng, and Cheng Guo, "A car-face region-based image retrieval method with attention of sift features," *Multimedia Tools and Applications (MTA), Springer*, pp. 1–20, 2016.

[14] Elad Hoffer and Nir Ailon, "Deep metric learning using triplet network," in *Similarity-Based Pattern Recognition*. 2015, pp. 84–92, Springer.

[15] Hao Ye, Zuxuan Wu, Rui-Wei Zhao, Xi Wang, Yu-Gang Jiang, and Xiangyang Xue, "Evaluating two-stream CNN for video classification," in *ACM Int. Conf. on Multimedia Retrieval*, 2015, ICMR, pp. 435–442.

[16] Dahjung Chung, Khalid Tahboub, and Edward J Delp, "A two stream siamese convolutional neural network for person re-identification," in *The IEEE international conference on computer vision (ICCV)*, 2017.

[17] Sergey Zagoruyko and Nikos Komodakis, "Learning to compare image patches via convolutional neural networks," in *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2015.

[18] D. C. Luvizon, B. T. Nassu, and R. Minetto, "A video-based system for vehicle speed measurement in urban roadways," *IEEE Trans. on Intelligent Transportation Systems (ITS)*, vol. PP, no. 99, pp. 1–12, 2016.

[19] Rodrigo Minetto, Nicolas Thome, Matthieu Cord, Neucimar J. Leite, and Jorge Stolfi, "T-HOG: An effective gradient-based descriptor for single line text regions," *Pattern Recognition (PR), Elsevier*, vol. 46, no. 3, pp. 1078–1090, 2013.

[20] Karen Simonyan and Andrew Zisserman, "Very deep convolutional networks for large-scale image recognition," *CoRR*, vol. abs/1409.1556, 2014.

[21] D. S. Bolme, J. R. Beveridge, B. A. Draper, and Y. M. Lui, "Visual object tracking using adaptive correlation filters," in *IEEE Computer Vision and Pattern Recognition (CVPR)*, June 2010, pp. 2544–2550.

[22] Xufeng Han, Thomas Leung, Yangqing Jia, Rahul Sukthankar, and Alexander. C. Berg, "Matchnet: Unifying feature and metric learning for patch-based matching," in *CVPR*, 2015.

[23] Jure Zbontar and Yann LeCun, "Stereo matching by training a convolutional neural network to compare image patches," *CoRR*, vol. abs/1510.05970, 2015.