

# Transfer Learning from Synthetic Data in the Camera Pose Estimation Problem

Jorge L. Charco<sup>1,2</sup>, Angel D. Sappa<sup>1,3</sup>, Boris X. Vintimilla<sup>1</sup> and Henry Velesaca<sup>1</sup>

<sup>1</sup>*ESPOL Polytechnic University, Escuela Superior Politécnica del Litoral, ESPOL. Campus Gustavo Galindo Km. 30.5 Vía Perimetral, P.O. Box 09-01-5863, Guayaquil, Ecuador*

<sup>2</sup>*University of Guayaquil, Delta and Kennedy Av., P.B. EC090514, Guayaquil, Ecuador*

<sup>3</sup>*Computer Vision Center, Edifici O, Campus UAB, 08193 Bellaterra, Barcelona, Spain  
{jlcharco, asappa, boris.vintimilla, hvelesac}@espol.edu.ec*

**Keywords:** Relative Camera Pose Estimation, Siamese Architecture, Synthetic Data, Deep Learning, Multi-View Environments, Extrinsic Camera Parameters.

**Abstract:** This paper presents a novel Siamese network architecture, as a variant of Resnet-50, to estimate the relative camera pose on multi-view environments. In order to improve the performance of the proposed model a transfer learning strategy, based on synthetic images obtained from a virtual-world, is considered. The transfer learning consist of first training the network using pairs of images from the virtual-world scenario considering different conditions (i.e., weather, illumination, objects, buildings, etc.); then, the learned weight of the network are transferred to the real case, where images from real-world scenarios are considered. Experimental results and comparisons with the state of the art show both, improvements on the relative pose estimation accuracy using the proposed model, as well as further improvements when the transfer learning strategy (synthetic-world data – transfer learning – real-world data) is considered to tackle the limitation on the training due to the reduced number of pairs of real-images on most of the public data sets.

## 1 INTRODUCTION

Automatic calibration of the camera extrinsic parameters is a challenging problem involved in several computer vision processes; applications such as driving assistance, human pose estimation, mobile robots, 3D object reconstruction, just to mention a few, would take advantage of knowing the relative pose between the camera and world reference system. During last decades different computer vision algorithms have been proposed for extrinsic camera parameters estimation (relative translation and rotation) (e.g., (Faugeras et al., 1992), (Hartley, 1994), (Sappa et al., 2006), (Liu et al., 2009), (Dornaika et al., 2011), (Schonberger and Frahm, 2016), (Iyer et al., 2018), (Lin et al., 2019)).

Classical approaches estimate the extrinsic camera parameters from feature points (e.g., SIFT (Lowe, 1999), SURF (Bay et al., 2006), BRIEF (Calonder et al., 2012)), which are detected and described in the pair of images used to estimate the pose (relative translation and rotation) between the cameras. These algorithms have low accuracy when they are unable to find enough common feature points to be matched be-

tween the images.

During last years, convolutional neural networks (CNNs) have been widely used for feature detection in tasks such as segmentation, images classification, super resolution and patten recognition, getting better results than state-of-art (e.g., (Kamnitsas et al., 2017), (Wang et al., 2016), (Rivadeneira et al., 2019)). Within this scheme, different CNN based camera calibrations, on single and multi-view environments, have been also proposed showing appealing results (e.g., (Shalnov and Konushin, 2017), (Charco et al., 2018)). The single view approaches capture real-world environments by using a single camera that is constantly moving around the scene that may contains dynamic objects (i.e., pedestrians, cars). The challenge with this approaches lies on the scenarios with moving objects, which change their position during the acquisition and could occlude different scene's regions (features) at consecutive frames. This problem is not present in multi-view approaches, due to the fact that the scene is simultaneously captured from different positions by different cameras, considering a minimum overlap of regions between the captured scenes. Fig. 1 shows an illustration of these two scenarios.

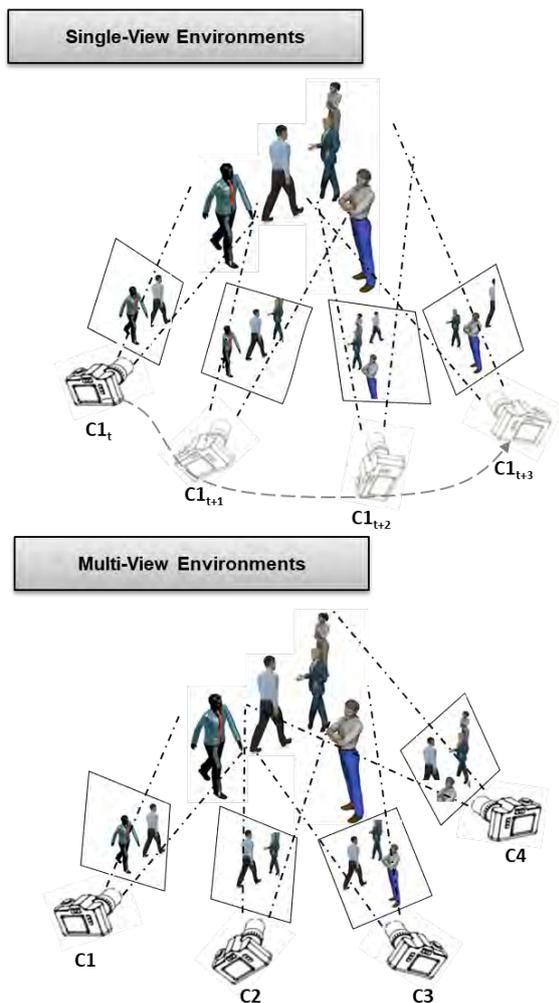


Figure 1: Illustration of single and multi-view environments.

Although appealing results are obtained with learning based approaches, they have as a main limitation the size of data set used for training the network, which becomes an important factor. In the particular case of extrinsic camera parameter estimation, there are several data sets devoted for this task. Actually, some of them have been proposed for applications such as mapping, tracking or relocalization (e.g., Cambridge (Kendall et al., 2015) and 7-Scene (Shotton et al., 2013)). Unfortunately, most of these data sets contain just a small number of images acquired by single view approaches. Trying to overcome this limitation, in (Aanæs et al., 2016), the authors have provided a multi-view data set, acquired using a robotic arm, called DTU-Robot. Main drawback with this data set lies on the fact that an indoor and small scenario, containing few objects (e.g.,

toys, woods, cans), has been considered; transferring knowledge from this scenario to real outdoor environments becomes a difficult task.

As mentioned above, the limited size of data sets is a common problem in learning-based computer vision solutions. In some cases this problem has been overcome by using virtual environments where data sets are acquired (e.g., pedestrian detection (Vázquez et al., 2014), image dehazing (Li et al., 2017), 3D object recognition (Jalal et al., 2019), optical flow estimation (Onkarappa and Sappa, 2015)). These virtual environments allow generating an almost unlimited set of synthetic images by rendered many times with different conditions and actors (i.e., weather, illumination, pedestrian, road, building, vehicles). Another appealing feature of synthetic data sets is related with the lack of manual ground truth annotations required in these scenarios. It allows saving time and resources needed by annotation tasks.

In the current work, a CNN architecture is proposed to estimate extrinsic camera parameters by using pairs of images acquired from the same scene and from different points of view at the same time. This architecture is firstly trained using synthetic images generated in a virtual world (i.e., a 3D representation of an urban environment containing a lot of buildings and different objects) and then its weights updated using a transfer learning strategy using real data. A Siamese network architecture is proposed as a variant of Resnet-50. The remainder of the paper is organized as follows. In Section 2 previous works on relative camera pose estimation are summarized; this is followed by Section 3, which presents the approach proposed for estimating extrinsic camera parameters and a detailed description of the used synthetic data sets. Experimental results are reported in Section 4 with the results obtained when the network is just trained with real-world data; comparisons with a previous approach are also presented. Finally, conclusion and future work are given in Section 5.

## 2 RELATED WORK

During the last years, CNNs models have been used in many computer vision tasks due to their capability to extract features improving state-of-art results. In that direction, some works have been proposed for camera pose estimation from a single view approach. The authors in (Kendall et al., 2015) have proposed a CNN architecture to regress the 6-DOF camera pose from a single RGB image, being robust to indoors and outdoors environments in real time, even difficult lighting, motion blur and different camera intrinsics pa-

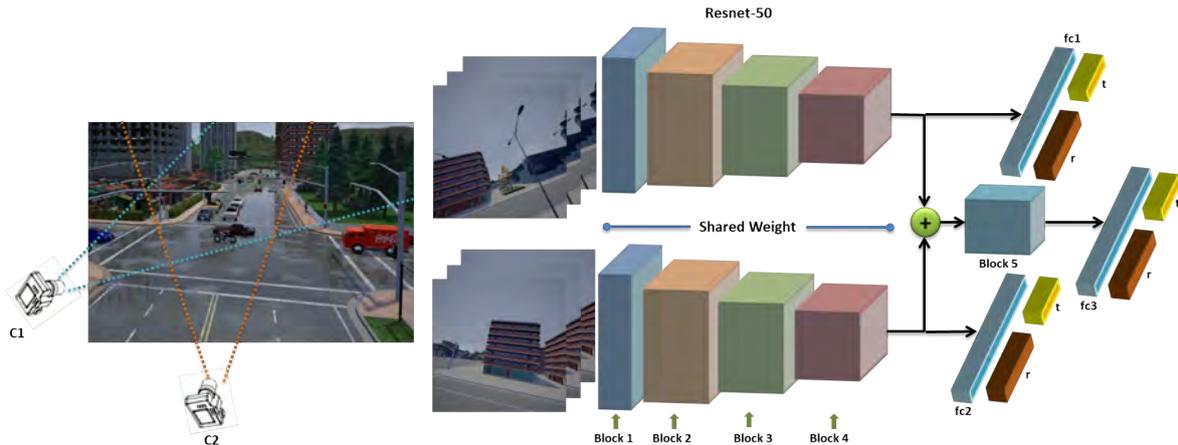


Figure 2: Siamese architecture feeds with two images of the same scene captured at the same time from different points of views. The regression part contains three fully-connected layers to estimate the extrinsic camera parameters.

rameters. The previous approach was updated with a similar architecture and a new loss function to learn camera pose in (Kendall and Cipolla, 2017). In (Shalnov and Konushin, 2017) the authors have proposed a novel method for camera pose estimation based on the scene’s prior-knowledge. The approach is trained on a synthetic scenario where the camera pose is estimated based on human body features. The trained network is then generalized to real environments. On the contrary to previous approaches, in (Iyer et al., 2018) the authors propose a self-supervised deep network to estimate the 6-DOF, rigid body transformation, between a 3D LiDAR and a 2D camera in real-time. The approach is then used to estimate the calibration parameters that maximize the geometric of the input images and point clouds.

On the other hand, just few works have been proposed to solve the camera pose estimation problem in multi-view environments using CNNs. In (Charco et al., 2018), the authors have proposed to use a Siamese CNN based on a modified AlexNet architecture with two identical branches and shared weights. The training process was performed from scratch with a set of pairs of images of the same scene simultaneously acquired from different points of view. The output of each branch is concatenated to two fully connected layers to estimate the relative camera pose (translation and rotation). Euclidean distance is used as a loss function. In (En et al., 2018) the authors have also proposed a Siamese Network with two branches regressing one pose per image. GoogLeNet architecture is used to extract features and the pose regressor contains two fully connected layers with ReLU activation. The quaternion is normalized during test time. Euclidean distance and weighting term ( $\beta$ ) are used to balance the error between translation and ro-

tation. The authors in (Lin et al., 2019) have presented an approach based on Recurrent Convolutional Neural Networks (RCNNs). They used the first four residual blocks of the ResNet-50. The output of each consecutive monocular image is concatenated to fed the last block of the ResNet-50. Two RCNNs are used, the first is fed by the concatenated output of two consecutive images to two Long Short-Term Memory (LSTM) to find the correlations among images. The second RCNN is fed from a monocular image to LSTM and its output is reshaped to a fully connected layer. Finally, both outputs are concatenated to obtain the translation and rotation estimation.

### 3 PROPOSED APPROACH

This section details the two main contributions of current work. Firstly, the Siamese network architecture proposed to estimate the relative pose between two cameras, which synchronously acquire images of the same scenario, is presented. The network takes as an input a pair of images of the same scenario, captured from different points of view (position and rotation) (see Fig. 2) and estimate the relative rotation and translation between the cameras. As mentioned above, the second contribution of current work lies on the strategy used to train the network. This strategy consists of training the proposed network using a synthetic data set (outdoor environments acquired by using CARLA simulator (Dosovitskiy et al., 2017)) and then transferring the knowledge of trained network in virtual environment to a real-world. The proposed transfer learning based strategy tackles the problem of having a large data set for the training process.

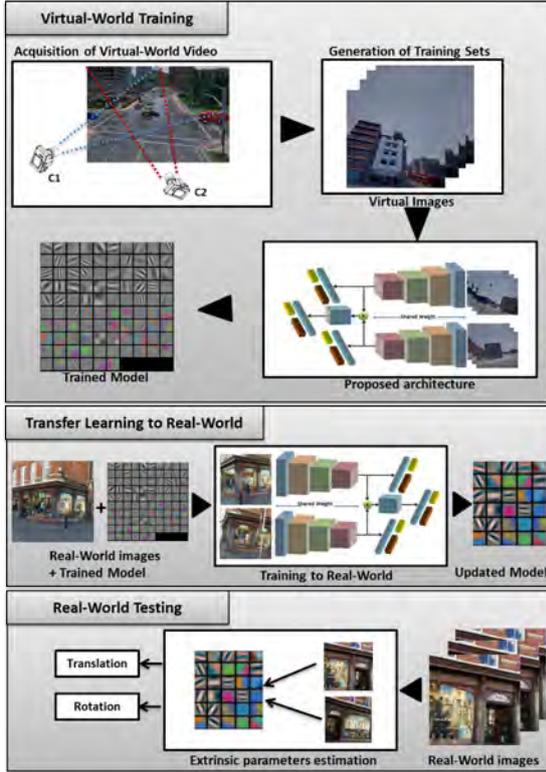


Figure 3: Training and testing processes.

### 3.1 Network Architecture

The proposed approach is a modified Resnet-50 (He et al., 2016), which contains two identical branches with shared weights up to the fourth residual block. The last residual (fifth block) is fed by concatenating the output of the fourth residual block from each branch. The architecture is composed with multiple residual units, bottleneck architecture consisting of convolutional layers, batch normalization, pooling and identity blocks. The standard residual block structure was modified by replacing RELUs with ELUs as activation function. According to (Clevert et al., 2015) ELU helps to speed up convergence and to reduce the bias shift neurons, avoiding the vanishing gradient. Furthermore, the last average pooling layer is replaced with a global average pooling layer. Two fully connected layers are added after the fourth residual blocks for each branch (left and right),  $fc1$  and  $fc2$ , and after the fifth residual block an additional fully connected layer is also added  $fc3$ . The **global pose** of each camera is predicted from the features extracted from the corresponding image, up to fourth residual block, which are used to feed the  $fc1$  or  $fc2$ . Each  $fc_i$  has a dimension of 1024 followed by two regressors, corresponding to the global trans-

lation ( $3 \times 1$ ) and rotation ( $4 \times 1$ ) respectively. Regarding the **relative camera pose**, the output of the last residual block (fifth block) is used to feed to  $fc3$  and subsequently the two regressors used to estimate the relative translation ( $3 \times 1$ ) and rotation ( $4 \times 1$ ) between the given pair of images.

Relative camera pose is represented by two vectors:  $\Delta p = [\hat{t}, \hat{r}]$ , where  $\hat{t}$  is a 3-dimensions vector that represents the translation and  $\hat{r}$  is a 4-dimensions vector that represents the rotation (i.e., a quaternion). As the image comes in the same context is reasonable to build one model to train both components (translation and rotation) at the same time. Normally, the Euclidean distance is used to calculate these components:

$$T_{Global}(I) = \|t - \hat{t}\|_{\gamma}, \quad (1)$$

$$R_{Global}(I) = \left\| r - \frac{\hat{r}}{\|\hat{r}\|} \right\|_{\gamma}, \quad (2)$$

where  $t$  represents the ground truth translation and  $\hat{t}$  denotes the prediction.  $r$  is the ground truth rotation and  $\hat{r}$  denotes the prediction of the quaternion values. The estimated rotation (i.e., the estimated quaternion values) is normalized to a unit length as  $\frac{\hat{r}}{\|\hat{r}\|}$ .  $\gamma$  is  $L2$  Euclidean norm.

The previous terms are merged according to a factor  $\beta$ , which is introduced due to the difference in scale between both components (translation and rotation) to balance the loss terms (Kendall et al., 2015). Hence, the general loss function, that includes both terms, is defined as:

$$Loss_{global}(I) = T_{Global} + \beta * R_{Global}, \quad (3)$$

setting the  $\beta$  parameter to the right value is a challenging task that depends on several factors related to the scene and cameras. In order to find the best solution (Kendall and Cipolla, 2017) propose to use two learnable variables called  $\hat{s}_x$  and  $\hat{s}_y$  that acts as weight to balance translation and rotation terms respectively—with a similar effect to  $\beta$ . Hence, in the current work, a modified loss function that uses  $\hat{s}_y$  as learnable variable is used:

$$Loss_{Global}(I) = T_{Global} + (\exp(\hat{s}_y) * R_{Global} + \hat{s}_y). \quad (4)$$

Each branch of the Siamese architecture estimates the global pose of the image. Additionally, by connecting these branches together through the fifth block the relative pose between the cameras is estimated; this relative pose is obtained as follow:

$$T_{Relative}(I) = \|t_{rel} - \hat{t}_{rel}\|_{\gamma}, \quad (5)$$

$$R_{Relative}(I) = \left\| r_{rel} - \frac{\hat{r}_{rel}}{\|\hat{r}_{rel}\|_\gamma} \right\|, \quad (6)$$

where  $T_{Relative}$  and  $R_{Relative}$  estimate the difference between the ground truth (i.e., the relative one) and the prediction of trained model ( $\hat{t}_{rel}$  and  $\hat{r}_{rel}$ ). As  $\hat{r}_{rel}$  is directly obtained from network it has to be normalized before. Equation (7) and Eq. (8) show how to obtain  $t_{rel}$  and  $r_{rel}$ .

$$t_{rel} = t_{C1} - t_{C2}, \quad (7)$$

$$r_{rel} = r_{C2}^* * r_{C1}, \quad (8)$$

where  $C_i$  corresponds to the pose parameters of the ( $i$ ) camera (i.e., rotation and translation) given by the CARLA simulator; these parameters are referred to a global reference system;  $r_{C2}^*$  is the conjugate quaternion of  $r_{C2}$ . Normalize the quaternion is an important process before using Eq. (8). Finally, the loss function used to obtain the relative pose of two images is defined as:

$$Loss_{Relative}(I) = T_{Rel} + (exp(\hat{s}_y) * R_{Rel} + \hat{s}_y). \quad (9)$$

Note that  $Loss_{Global}$  in Eq. (4) and  $Loss_{Relative}$  in Eq. (9) are applied for different purposes. The first one is used to predict global pose through each branch of the trained model. While the second predicts the relative pose by concatenating the Siamese Network (see Fig. 2). The proposed approach was jointly trained with Global and Relative Loss, as shown in Eq. (10):

$$L = Loss_{Global} + Loss_{Relative}. \quad (10)$$

### 3.2 Synthetic Data Set

This section presents the steps followed for the synthetic data set generation. Two open-source software tools were used: CARLA Simulator (Dosovitskiy et al., 2017) and OpenMVG (Moulon et al., 2016). This first one, CARLA, is used for generating synthetic images from a virtual-world. It has been developed from the ground up to support the development, training, and validation of autonomous urban driving systems. The simulation platform supports the flexible specification of sensor suites and environmental conditions (Dosovitskiy et al., 2017). The second open-source software (OpenMVG) has been used to estimate the overlap between a given pair of images. It is designed to provide easy access to the classical problem solvers in multiple view geometry and solve them accurately.

A workstation with a Titan XP GPU was used for server execution due to the large amount of processing demanded by the CARLA simulator. The attributes and values used to configure the cameras in CARLA simulator are shown in Table 1. The cameras were attached to moving platform to allow navigation through the virtual-world. Fifteen different types of weathers are used to generate the set of pairs of images. All selected weather and their respective parameter settings are shown in Table 2.

Table 1: Attributes and values for the two virtual cameras in CARLA.

Attribute	Value
image_size_x	448
image_size_y	448
sensor_tick	1.0
fov	100.0



Figure 4: Trajectory followed by the camera in CARLA Simulator.

The trajectory followed by the cameras is shown in Fig 4. It starts from the green point and moves to left following the arrows until it ends at the same green point where it starts. The pair of cameras follow the aforementioned path while their relative pose (relative position and orientation between them) randomly changes at each time in the range  $x = [-0.5, 0.5]$ ,  $y = [-0.5, 0.5]$ , and  $z = [-0.5, 0.5]$ , while their relative orientation changes in the range  $roll = [0, 12]$ ,  $pitch = [-5, 5]$ , and  $yaw = [-5, 5]$  degrees, these values were defined according to the scene and building characteristics. Values outside these ranges generate a pair of images with little overlap.

OpenMVG is used to obtain the list of pairs of images with an overlap higher than a give threshold. The configuration options of OpenMVG has been set as presented in Table 3. After running the OpenMVG script, the pairs of images with an overlap higher than the given threshold has been obtained, according to the parameters mentioned above. The images are captured at the size of  $448 \times 448$ . The ground truth for the

Table 2: Attributes and values for the selected weathers.

Weather	Cloudyness	Precipitation	Precipitation deposits	Wind intensity	Sun azimuth (ang.)	Sun altitude (ang.)
Custom weather	0	0	0	0.00	-90	60
Clear noon	15	0	0	0.35	0	75
Cloudy noon	80	0	0	0.35	0	75
Wet noon	20	0	50	0.35	0	75
Wet cloudy noon	80	0	50	0.35	0	75
Mid rainy noon	80	30	50	0.40	0	75
Hard rainy noon	90	60	100	1.00	0	75
Soft rain noon	70	15	50	0.35	0	75
Clear sunset	15	0	0	0.35	0	15
Cloudy sunset	80	0	0	0.35	0	15
Wet sunset	20	0	50	0.35	0	15
Wet cloudy sunset	90	0	50	0.35	0	15
Mid rain sunset	80	30	50	0.40	0	15
Hard rain sunset	80	60	100	1.00	0	15
Soft rain sunset	90	15	50	0.35	0	15

global camera pose has been obtained from CARLA simulator, while the relative camera pose is computed from Eq. (7) and Eq. (8). Additionally, each pair of images, and their ground truth camera pose (global pose and relative pose), has been saved into TFRecord file of tensorflow, for a better memory management during training. A summary of the data set obtained at the different weather conditions is presented in Table 4.

Table 3: Options available for OpenMVG script.

Order	Option	Used
1	Intrinsics analysis	Yes
2	Compute features	Yes
3	Compute matches	Yes
4	Do Incremental/Sequential reconstruction	No
5	Colorize Structure	No
6	Structure from Known Poses (robust triangulation)	No

## 4 EXPERIMENTS RESULTS

As mentioned above, this paper has two main contributions. On the one hand, a new CNN based architecture is proposed for extrinsic camera parameter estimation; on the other hand, trying to overcome limitations related with the reduced amount of data provided in most of public data sets a transfer learning strategy is proposed. This strategy is based on the usage of synthetic data generated using CARLA simulator (Dosovitskiy et al., 2017). Hence, this section first presents details on the data set generation; then, experimental results by training the proposed architecture with the ShopFacade and OldHospital of Cam-

Table 4: Data set and weather conditions.

Weather	Pair images
Custom weather	2545
Clear noon	1780
Cloudy noon	2002
Wet noon	1703
Wet cloudy noon	1232
Mid rainy noon	1641
Hard rain noon	1323
Soft rain noon	1213
Clear sunset	1136
Cloudy sunset	1144
Wet sunset	1450
Wet cloudy sunset	1512
Mid rain sunset	1672
Hard rain sunset	1685
Soft rain sunset	1346
<b>Total</b>	<b>23384</b>

bridge data set (Kendall et al., 2015) are depicted; and finally, results obtained when the proposed transfer learning strategy is used (i.e., first training the network on the synthetic data and then updating network weights by keep training it with real data).

The proposed approach was implemented with the TensorFlow library (Abadi et al., 2016) and trained with NVIDIA Titan XP GPU and Intel Core I9 3.3GHz CPU. Adam optimizer (Kingma and Ba, 2014) is used to train the network with a learning rate of  $10^{-4}$  and batch size of 32. The  $\hat{s}_y$  is initialized with -6.0 in all the experiments.

### 4.1 Training on Real Data

Using the strategy of transfer learning, all layers were initialized up to the fourth residual block with the weights of Resnet-50 pretrained on ImageNet and the normal distribution initialization was used for the remaining layers. The network architecture was trained

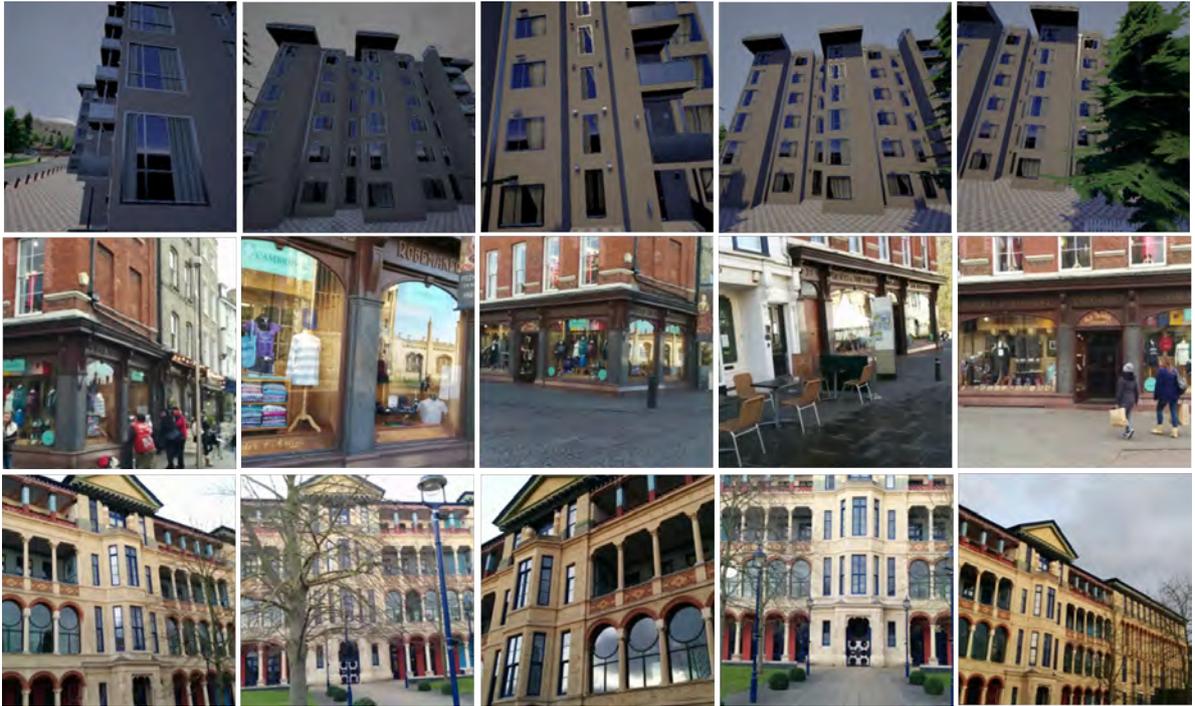


Figure 5: (1st row) Images from the synthetic image data set generated by the CARLA simulator (they were captured at different weather and lighting conditions. (2nd and 3rd row) Real-world images, ShopFacade and OldHospital of Cambridge data set respectively, used to evaluate the trained model.

on ShopFacade and OldHospital of Cambridge data set. As pre-processing data set, the images were resized to 224 pixels along the shorter side; then, the mean value was computed and subtracted from the images. For the training process, random crops of  $224 \times 224$  pixels have been computed on the OldHospital data set; this results in a set of 5900 pairs of images, which were used to feed to the network. The same process has been performed at 1300 pairs of images of ShopFacade data set. Both data set were used to train the network until 500 epochs, which approximately took 7 hours and 3 hours respectively.

The pre-processing mentioned above has been also used during the evaluation phase. In the evaluation a set of 2100 pairs of images from the OldHospital and a set of 250 pairs of images from ShopFacade have been considered. On the contrary to the training stage, in this case a central crop is used instead of a random crop. Since the relationship between the pair of images is required for the training stage, both the relative and the absolute pose were estimated.

## 4.2 Transfer Learning Strategy

This section presents details on the strategy followed for transferring the parameters learned in a synthetic

data set to real images. The proposed architecture was initialized as presented in Section 4.1. In this case the training process was performed on the synthetic data set described in Section 3.2. Since the size of synthetic images is  $448 \times 448$  pixels, a resize up to  $224 \times 224$  pixels is performed; like in the previous case, the average mean value is estimated and subtracted from each image. The synthetic data set contains 23384 pairs of images, which were used to feed to the architecture. The training stage took about 5 hours with 300 epochs. After the training process from the synthetic images, learned weights were used to initialize all layers of proposed architecture, which is re-trained and refined in an end-to-end way with real-world images.

## 4.3 Results

Experimental results obtained with the proposed network and training strategy are presented. Additionally, the obtained results are compared with a state-of-the-art CNN-based method Pose-MV (Charco et al., 2018) on ShopFacade and OldHospital of Cambridge data sets. Average median error on rotation and translation for both data sets are depicted on Table 5. Angular error and Euclidean distance error are used to

Table 5: Comparison of average median Relative Pose Errors (extrinsic parameters) of RelPoseTL with existing CNN models on ShopFacade and OldHospital of Cambridge data set.

Scene / Models	State-of-the-art	RelPoseTL (Ours experiments)	
	Pose-MV	Real data	Sythetic data (Transfer Learning)
ShopFacade	1.126m, 6.021°	1,002m, 3.655°	<b>0.834m, 3.182°</b>
OldHospital	5.849m, 7.546°	3.792m, 2.721°	<b>3.705m, 2.583°</b>
Average	3.487m, 6.783°	2.397m, 3.188°	<b>2.269m, 2.882°</b>

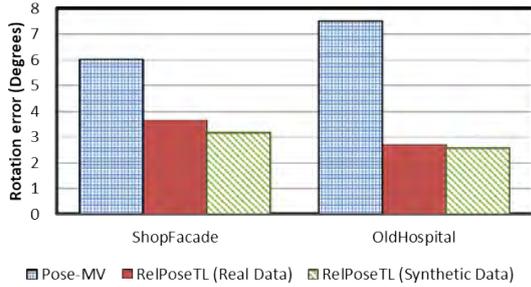


Figure 6: Comparisons on **rotation** error between Pose-MV and the proposed RealPoseTL on ShopFacade and OldHospital of Cambridge data set.

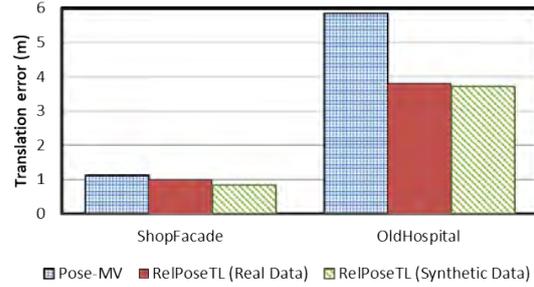


Figure 7: Comparisons on **translation** error between Pose-MV and the proposed RealPoseTL on ShopFacade and OldHospital of Cambridge data set.

evaluate the performance of the proposed approach. The first one is used to compute the rotation error between the obtained result and the given ground truth value using the 4-dimensional vector (quaternion). The second one is used to measure the distance error between the estimated translation and the corresponding ground truth (3-dimensional translation vector). The proposed approach obtains more accurate results on both translation and rotation in both data sets. The average median translation error improves by 32% the results obtained by RelPoseTL(Real Data) with respect to the results obtained with Pose-MV (both approaches trained with the same real data); this improvements reaches up to 35% when the proposed transfer learning strategy is considered. With respect to average median rotation error, the results obtained with the proposed RelPoseTL (Real Data) improves by 53% the results obtained with Pose-MV; this improvements reaches up to the 58% when the proposed transfer learning strategy from synthetic data is considered. Figures 6 and 7 illustrate the comparisons between the proposed approach and Pose-MV on ShopFacade and OldHospital of Cambridge data set. The results with large translation and rotation errors correspond to pairs of images with large relative motions between them (see Fig. 8).

## 5 CONCLUSION

This paper addresses the challenging problem of estimating the relative camera pose from two different images of the same scenario, acquired from two different points of view at the same time. A novel deep learning based architecture is proposed to accurately estimate the relative rotation and translation between the two cameras. Experimental results and comparisons are provided showing improvements on the obtained results. As a second contribution of this paper a training strategy based on transfer learning from Synthetic data is proposed. This strategy is motivated by the reduced amount of images on the data sets provided to train the network. The manuscript shows how features extracted from large amounts of synthetic images can help the estimation of relative camera pose in real-world images. The proposed architecture has been trained both, by only using real images and by using the proposed transfer learning strategy. Experimental results show that the proposed transfer learning approach helps to reduce error in the obtained results (relative rotation and translation). Future work will be focused on extending the usage the synthetic images data set, by increasing the data set with others outdoor multi-view environments, including different virtual simulators. The goal of this future work is to reach the state-of-the-art results without using a transfer learning strategy. Additionally, different

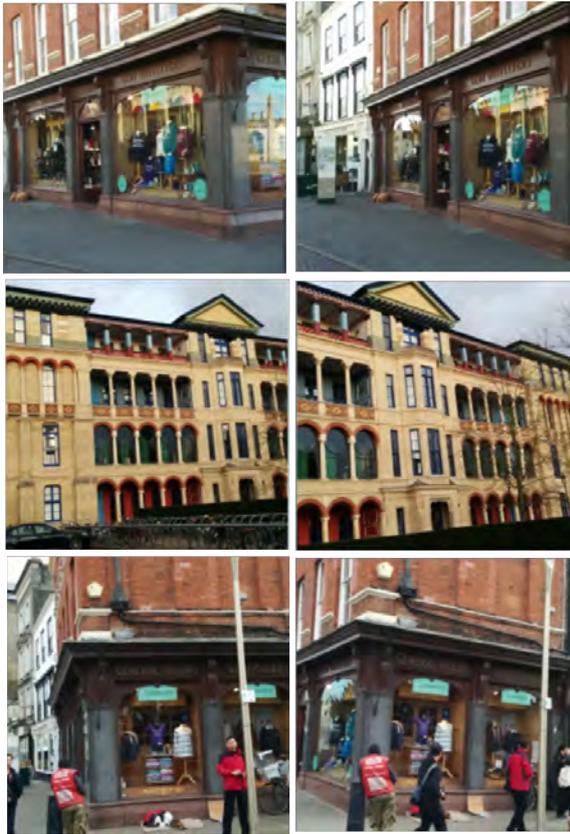


Figure 8: Challenging scenarios, pairs of images from points of view with large relative rotation and translation.

models based on CNNs will be considered.

## ACKNOWLEDGEMENTS

This work has been partially supported by the ESPOL project EPASI (CIDIS-01-2018) and Aplicaciones TICs para Ciudades Inteligentes (FIEC-16-2018); the Spanish Government under Project TIN2017-89723-P; and the “CERCA Programme / Generalitat de Catalunya”. The authors gratefully acknowledge the support of the CYTED Network: “Ibero-American Thematic Network on ICT Applications for Smart Cities” (REF-518RT0559) and the NVIDIA Corporation for the donation of the Titan Xp GPU used for this research. The first author has been supported by Ecuador government under a SENESCYT scholarship contract.

## REFERENCES

Aanaes, H., Jensen, R. R., Vogiatzis, G., Tola, E., and Dahl, A. B. (2016). Large-scale data for multiple-view

stereopsis. *International Journal of Computer Vision*, pages 1–16.

Abadi, M., Barham, P., Chen, J., Chen, Z., Davis, A., Dean, J., Devin, M., Ghemawat, S., Irving, G., Isard, M., et al. (2016). Tensorflow: A system for large-scale machine learning. In *12th {USENIX} Symposium on Operating Systems Design and Implementation ({OSDI} 16)*, pages 265–283.

Bay, H., Tuytelaars, T., and Gool, L. J. V. (2006). SURF: Speeded Up Robust Features. In *Proceedings of the 9th European Conference on Computer Vision, Graz, Austria, May 7-13*, pages 404–417.

Calonder, M., Lepetit, V., Özuysal, M., Trzcinski, T., Strecha, C., and Fua, P. (2012). BRIEF: Computing a local binary descriptor very fast. *IEEE Trans. Pattern Anal. Mach. Intell.*, 34(7):1281–1298.

Charco, J. L., Vintimilla, B. X., and Sappa, A. D. (2018). Deep learning based camera pose estimation in multi-view environment. In *2018 14th International Conference on Signal-Image Technology & Internet-Based Systems (SITIS)*, pages 224–228. IEEE.

Clevert, D.-A., Unterthiner, T., and Hochreiter, S. (2015). Fast and accurate deep network learning by exponential linear units (elus). *arXiv preprint arXiv:1511.07289*.

Dornaika, F., Álvarez, J. M., Sappa, A. D., and López, A. M. (2011). A new framework for stereo sensor pose through road segmentation and registration. *IEEE Transactions on Intelligent Transportation Systems*, 12(4):954–966.

Dosovitskiy, A., Ros, G., Codevilla, F., Lopez, A., and Koltun, V. (2017). CARLA: An open urban driving simulator. In *Proceedings of the 1st Annual Conference on Robot Learning*, pages 1–16.

En, S., Lechervy, A., and Jurie, F. (2018). Rpnnet: an end-to-end network for relative camera pose estimation. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 0–0.

Faugeras, O. D., Luong, Q. T., and Maybank, S. J. (1992). Camera self-calibration: Theory and experiments. In Sandini, G., editor, *Computer Vision — ECCV’92*, pages 321–334, Berlin, Heidelberg. Springer Berlin Heidelberg.

Hartley, R. I. (1994). Self-calibration from multiple views with a rotating camera. In *European Conference on Computer Vision*, pages 471–478. Springer.

He, K., Zhang, X., Ren, S., and Sun, J. (2016). Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778.

Iyer, G., Ram, R. K., Murthy, J. K., and Krishna, K. M. (2018). Calibnet: Geometrically supervised extrinsic calibration using 3d spatial transformer networks. In *2018 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 1110–1117. IEEE.

Jalal, M., Spjut, J., Boudaoud, B., and Betke, M. (2019). Sidod: A synthetic image dataset for 3d object pose recognition with distractors. In *Proceedings of the*

- IEEE Conference on Computer Vision and Pattern Recognition Workshops*, pages 0–0.
- Kamnitsas, K., Ledig, C., Newcombe, V. F., Simpson, J. P., Kane, A. D., Menon, D. K., Rueckert, D., and Glocker, B. (2017). Efficient multi-scale 3d cnn with fully connected crf for accurate brain lesion segmentation. *Medical image analysis*, 36:61–78.
- Kendall, A. and Cipolla, R. (2017). Geometric loss functions for camera pose regression with deep learning. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 5974–5983.
- Kendall, A., Grimes, M., and Cipolla, R. (2015). Posenet: A convolutional network for real-time 6-dof camera relocalization. In *Proceedings of the IEEE international conference on computer vision*, pages 2938–2946.
- Kingma, D. P. and Ba, J. (2014). Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*.
- Li, K., Li, Y., You, S., and Barnes, N. (2017). Photo-realistic simulation of road scene for data-driven methods in bad weather. In *The IEEE International Conference on Computer Vision (ICCV) Workshops*.
- Lin, Y., Liu, Z., Huang, J., Wang, C., Du, G., Bai, J., and Lian, S. (2019). Deep global-relative networks for end-to-end 6-dof visual localization and odometry. In *Pacific Rim International Conference on Artificial Intelligence*, pages 454–467. Springer.
- Liu, R., Zhang, H., Liu, M., Xia, X., and Hu, T. (2009). Stereo cameras self-calibration based on sift. In *2009 International Conference on Measuring Technology and Mechatronics Automation*, volume 1, pages 352–355.
- Lowe, D. G. (1999). Object recognition from local scale-invariant features. In *Proceedings of the IEEE International Conference on Computer Vision, Kerkyra, Greece, September 20-27*, pages 1150–1157.
- Moulon, P., Monasse, P., Marlet, R., and Others (2016). Openmvg. an open multiple view geometry library. <https://github.com/openMVG/openMVG>.
- Onkarappa, N. and Sappa, A. D. (2015). Synthetic sequences and ground-truth flow field generation for algorithm validation. *Multimedia Tools and Applications*, 74(9):3121–3135.
- Rivadeneira, R. E., Suárez, P. L., Sappa, A. D., and Vintimilla, B. X. (2019). Thermal image superresolution through deep convolutional neural network. In *International Conference on Image Analysis and Recognition*, pages 417–426. Springer.
- Sappa, A., Gerónimo, D., Dornaika, F., and López, A. (2006). On-board camera extrinsic parameter estimation. *Electronics Letters*, 42(13):745–747.
- Schonberger, J. L. and Frahm, J.-M. (2016). Structure-from-motion revisited. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 4104–4113.
- Shalnov, E. and Konushin, A. (2017). Convolutional neural network for camera pose estimation from object detections. *International Archives of the Photogrammetry, Remote Sensing & Spatial Information Sciences*, 42.
- Shotton, J., Glocker, B., Zach, C., Izadi, S., Criminisi, A., and Fitzgibbon, A. (2013). Scene coordinate regression forests for camera relocalization in rgb-d images. In *2013 IEEE Conference on Computer Vision and Pattern Recognition*, pages 2930–2937.
- Vázquez, D., López, A. M., Marín, J., Ponsa, D., and Gerónimo, D. (2014). Virtual and real world adaptation for pedestrian detection. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 36(4):797–809.
- Wang, J., Yang, Y., Mao, J., Huang, Z., Huang, C., and Xu, W. (2016). Cnn-rnn: A unified framework for multi-label image classification. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2285–2294.