

# Visual Localization Under Appearance Change: A Filtering Approach

Anh-Dzung Doan, Yasir Latif, Tat-Jun Chin, Yu Liu, Shin Fang Ch'ng, Thanh-Toan Do, and Ian Reid

**Abstract**—A major focus of current research on place recognition is visual localization for autonomous driving. In this scenario, as cameras will be operating continuously, it is realistic to expect videos as an input to visual localization algorithms, as opposed to the single-image querying approach used in other place recognition works. In this paper, we show that exploiting temporal continuity in the testing sequence significantly improves visual localization - qualitatively and quantitatively. Although intuitive, this idea has not been fully explored in recent works. Our main contribution is a novel Monte Carlo-based visual localization technique that can efficiently reason over the image sequence. Also, we propose an image retrieval pipeline which relies on local features and an encoding technique to represent an image as a single vector. The experimental results show that our proposed method achieves better results than state-of-the-art approaches for the task on visual localization under significant appearance change. Our synthetic dataset and source code are publicly made available<sup>1 2</sup>

**Index Terms**—place recognition, autonomous driving, Monte Carlo localization, image retrieval

## I. INTRODUCTION

To carry out higher level tasks such as planning and navigation, a robot needs to maintain, at all times, an accurate estimate of its position and orientation with respect to the environment. When the robot uses an existing map to infer its 6 degree of freedom (DoF) pose, the problem is termed *localization*. In the case when the map information is appearance (images) associated with different parts of the map, the problem is that of *visual localization* (VL). Image-based localization methods normally assume that the appearance remains unchanged from when the map is generated to the present time when the robot needs to localize itself. However, as the operational time span of the robot increases, the appearance of the environment inevitably changes. This poses a great challenge for visual localization methods as the underlying assumption of static appearance is violated due to continual changes of environment, e.g., weather condition, time of day, construction sites, updating of façades and billboards, etc.

One approach towards dealing with appearance change is to observe as many variations as possible of each location and carry a representation that can model them [1]. However, the bottleneck in this case is the significant amount of time needed

to capture sufficient naturally occurring appearance variations. For example, recent benchmarks [2]–[4] have taken years to collect data with sufficient natural appearance variation.

An orthogonal approach to address appearance change, pioneered by SeqSLAM [5], is to consider sequence-to-sequence matching instead of matching a single query image to previously observed images. SeqSLAM showed that sequence-based matching is more resilient to appearance changes but comes with its own shortcoming (e.g., sensitivity to view-point changes and differences in sequence velocities [6]). While recent learning-based VL algorithms [7]–[10] have focused on the case of a single query image, the camera, in robotics setting, is operating continuously once the vehicle is moving and hence, it is realistic to expect a video (sequence of images) as an input to VL methods.

Another paradigm to overcome this challenge is to extract “fingerprints” which are robust against appearance changes [11] [12]. However, this approach has been only shown its efficiency in “natural” changes, e.g., weather conditions or time of day.

**Our contribution:** This work addresses the problem of metric visual localization (i.e., we recover the 6 DoF pose of the camera) under appearance changes, in the context of autonomous vehicles. To effectively reason over image sequences, we leverage Monte Carlo principle that approximates the probability distribution of 6 DoF camera poses incrementally. We show that for the case of driving on urban streets, a simple motion model is enough to successfully track the vehicle’s pose. For the observation model, we propose a novel *observation encoder* based on image retrieval which generates a fixed (low) dimensional vector representation for each image. We show experimentally that such a representation is more resilient to appearance change along with minor changes in the viewpoint, while being compact in size.

We validate the proposed localization method on synthetic and real datasets [4] to show that it outperforms state-of-the-art approaches in the task of 6 DoF visual localization in the large-scale environment with significant appearance changes, i.e., the traversal distance of the vehicle is about 10 km.

## II. RELATED WORKS

Visual localization is a very well-studied area and various methods have been proposed to address the problem in different ways. For the scope of our work, we limit the survey to methods that can perform metric localization, i.e., they recover the 6 DoF pose of the camera based on image information.

Anh-Dzung Doan, Yasir Latif, Tat-Jun Chin, Yu Liu, Shin-Fang Ch'ng and Ian Reid are with School of Computer Science, The University of Adelaide. Email: {dung.doan, yasir.latif, tat-jun.chin, yu.liu04, shinfang.chng, ian.reid}@adelaide.edu.au

Thanh-Toan Do is with Department of Computer Science, University of Liverpool. Email: thanh-toan.do@liverpool.ac.uk

<sup>1</sup><http://tiny.cc/jd73bz>

<sup>2</sup><https://github.com/dadung/Visual-Localization-Filtering>

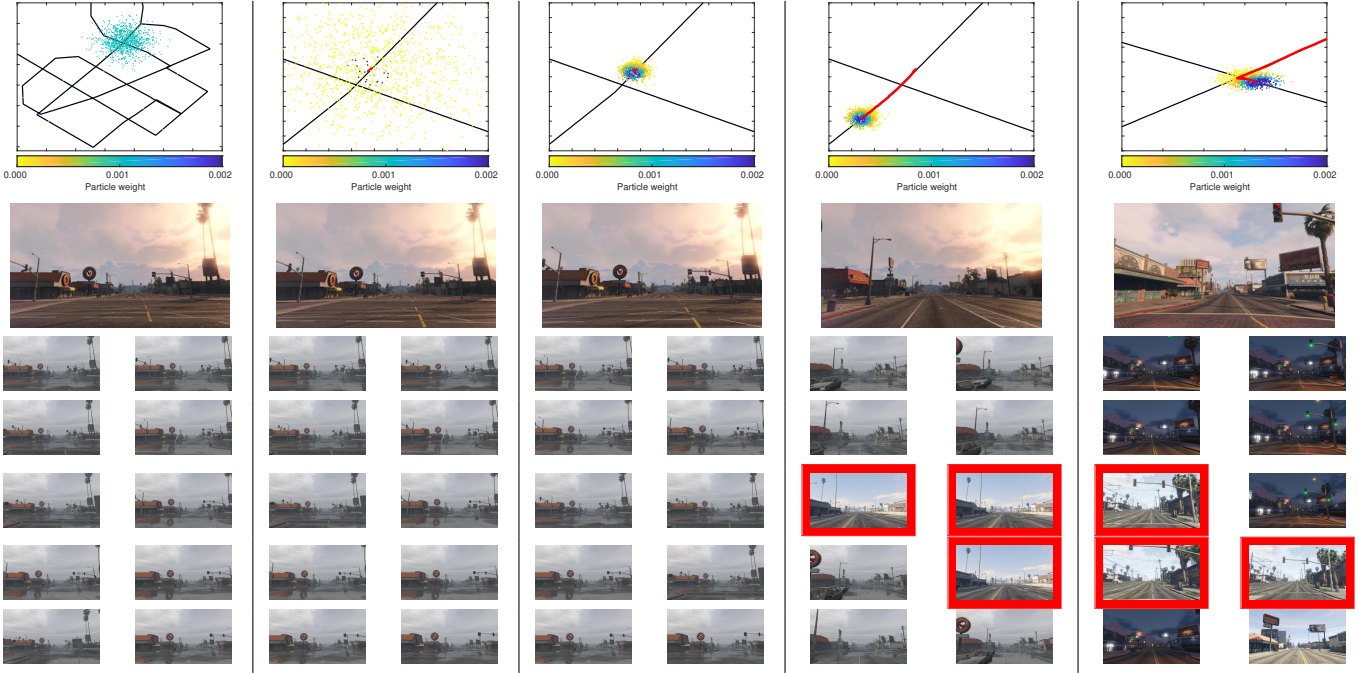


Fig. 1: An example of our proposed method. 1st and 2nd rows respectively are the particle distribution and query image, from 3rd to 7th rows are the retrieval results. Red bounding boxes indicate mistakes. The color bar describes the particle weights. Red lines are the camera pose predicted by our method. At the first frame, we randomly generate 1,000 particles with weight 0.001, then particles which are inconsistent with measurements are vanished. Eventually, the ambiguity is resolved, and the particles track the movement of the vehicle. It is also worth mentioning that the query and retrieved images are from different time and weather conditions. In particular, the query images are from 6:26pm-cloudy, and the correct retrieved images are from 11:24am-rainy, 9:15pm-clear, and 1:28pm-sunny conditions.

These algorithms can be categorized into local feature-based methods and learning-based methods.

Broadly speaking, local feature-based methods estimate the pose of the query image from 2D-3D correspondences by solving the PnP problem [13] or an equivalent step [2]. The main difficulty for this approach is establishing and validating the 2D-3D correspondences for large (city-scale) maps. Many methods attempt to improve this step by using geometric constraints [7], [14], semantic information in the images [15], and learning-based matching [9].

Learning-based methods perform 6-DoF pose estimation as an image-to-SE(3) regression problem. PoseNet [8] uses a convolutional neural network (CNN) to learn the regression mapping. LSTM-Pose [16] argues that the usage of a fully connected layer in PoseNet possibly leads to the overfitting, despite the usage of dropout. Therefore, they use Long-Short Term Memory (LSTM) units. Other variants of PoseNet include: uncertainty in the pose prediction [17] and geometric constraints [10], [18].

Our proposed method is different from above methods, in that we seek not only the current 6 DoF pose of the camera as a point estimate but as a distribution over the space of all possible locations in the map. Therefore, we employ a probabilistic filtering framework, namely Monte Carlo localization, that takes into account the constraints on motion of

the vehicle and the temporal nature of the localization task. To deal with appearance changes, we propose a novel observation encoder based on image retrieval which relies on dense local features and encoding technique to ensure the robustness against environmental changes. While, previous works [19] [20] [21] have combined Monte Carlo localization with image retrieval for localizing a robot in an indoor environment, they only consider the case of a robot moving on plane (3 DoF) while we aim to recover the complete 6 DoF pose of the robot. Also, while these works [19] [20] [21] only test their method in the indoor environment, in which the appearance changes are insignificant; we show that our method can perform robustly in a large-scale outdoor environment with significant appearance changes, i.e., the car traverses over 10 km in its run.

### III. MONTE CARLO-BASED VISUAL LOCALIZATION

Let the 6 DoF camera pose be given by:

$$s_t = [r_t, \Omega_t]^T$$

where,  $r_t$  and  $\Omega_t$  represent the 3D position and Euler orientation respectively at time  $t$ . Given motion  $u_{1:t}$  and noisy measurement  $z_{1:t}$  up to time  $t$ , we would like to estimate probability distribution  $p(s_t|u_{1:t}, z_{1:t})$ . However,  $p(s_t|u_{1:t}, z_{1:t})$  can be an arbitrary distribution, thus we leverage Monte Carlo principle to address this issue.

The idea of Monte Carlo-based visual localization is to represent the probability distribution  $p(s_t|u_{1:t}, z_{1:t})$  with a set of  $N$  particles. Each particle maintain an estimate of the state at time  $t$ :

$$\mathcal{S}_t = \{s_t^{[1]}, s_t^{[2]}, \dots, s_t^{[N]}\}$$

with a set of corresponding weights:

$$\mathcal{W}_t = \{w_t^{[1]}, w_t^{[2]}, \dots, w_t^{[N]}\}$$

The states of particles and their corresponding weights are respectively updated according to the motion  $u_t$  (Section III-A) and noisy measurement  $z_t$  (Section III-B) at time  $t$ . In the following sections, we justify the use of a simple motion model and give details about the novel observation encoder.

#### A. Motion model

In the autonomous driving scenario when navigating the roads of an urban area<sup>3</sup>, the motion of the car is fairly restricted, i.e., it largely stays on a road network on an approximately 2D plane<sup>4</sup> [23]. While the road surface still allows significant scope for movement (cf. Fig. 2), relative to the size of the map, the motion of the car is rather limited<sup>5</sup>. This is echoed by the recent work of [2], who observed that there is “*lower variation in viewpoints as the car follows the same road*”. Hence, a Monte Carlo scheme with a simple motion model suffices to track the 6 DoF rigid motion of the vehicle. In many cities, the road networks are complex Euclidean graphs. In fact, it is well known that using (visual) odometry alone, it is possible to accurately track a car on a 2D road map [24]<sup>6</sup>. More fundamentally, this suggests that temporal continuity in the testing sequence (which is fully exploited by our method) strongly benefits VL.

Mathematically, for each particle, we model the noisy motion consisting of linear velocity  $v_t^{[i]}$  and angular velocity  $\psi_t^{[i]}$  as the following:

$$\begin{aligned} v_t^{[i]} &\sim \mathcal{N}(\mu_v, \Sigma_v) \\ \psi_t^{[i]} &\sim \mathcal{N}(\mu_\psi, \Sigma_\psi) \end{aligned} \quad (1)$$

where,  $\mu_v$  and  $\mu_\psi$  respectively represent the linear and angular velocities. The accelerations are modeled by the noise covariance matrices  $\Sigma_v$  and  $\Sigma_\psi$ . For each particle, their motion in each time step is given by:

$$u_t^{[i]} = [v_t^{[i]}, \psi_t^{[i]}]^T$$

In practice, the  $\mu_v$ ,  $\mu_\psi$ ,  $\Sigma_v$ , and  $\Sigma_\psi$  can be either manually tuned, or estimated from training data [25].

<sup>3</sup>In more “localized” operations such as parking, where highly accurate 6-DoF estimation is required, it is probably better to rely on the INS.

<sup>4</sup>More fundamentally, the car is a non-holonomic system [22].

<sup>5</sup>On uneven or hilly roads, accelerometers can be used to estimate the vertical motion, hence, VL can focus on map-scale navigation.

<sup>6</sup>The method of [24] will give ambiguous results on non-informative trajectories, e.g., largely straight routes. Hence, VL is still crucial.

While 3D positions can be easily updated by using simple additions, we convert two Euler angles to the Direction Cosine Matrix (DCM) [26], multiply two matrices and convert the result back to the Euler representation, to stay on the 3D manifold of valid rotations. Let  $\varphi(\cdot)$  be a function that maps an Euler representation to DCM and  $\varphi^{-1}(\cdot)$  is its inverse mapping. Our motion model for each particle is then given by:

$$s_t^{[i]} = \begin{bmatrix} r_{t-1}^{[i]} + v_t^{[i]} \\ \varphi^{-1}(\varphi(\psi_t^{[i]}) \cdot \varphi(\Omega_{t-1}^{[i]})) \end{bmatrix} \quad (2)$$

The experimental results will show that our motion model can properly handle the temporal evolution of camera pose in an image sequence. In case there is mismatch between the actual motion and the one predicted by the motion model, such as during emergency maneuvers, the discrepancy would be reflected in the enlarged covariance estimate and resolved once motion returns to within normal bounds.

#### B. Observation encoder based on image retrieval

To make the image representation robust against appearance change, we seek for every image  $I$  a nonlinear function  $\tau(I)$  that maps  $I$  to a vector in a fixed dimensional space. To do so, we first densely compute SIFT features:  $\{x_i \in \mathbb{R}^d \mid i = 1, \dots, M\}$  over the image, followed by RootSIFT normalization [27]:

- i) L1 normalize every SIFT vector  $x_i = x_i / \|x_i\|_1$
- ii) square root each element  $x_i = \sqrt{x_i}$ .

where, RootSIFT normalization makes the Euclidean distance calculation among SIFT features equivalent to computing the Hellinger kernel.

Subsequently, we employ the embedding function VLAD [28] to embed SIFT features into a higher dimensional vector space. In particular, given a vocabulary learned by K-means:  $\mathcal{C} = \{c_k \in \mathbb{R}^d \mid k = 1, \dots, K\}$ , every SIFT feature is embedded as follows:

$$\phi_{VLAD}(x_i) = [\dots, 0, x_i - c_j, 0, \dots] \in \mathbb{R}^D$$

where  $c_j$  is the nearest visual word to  $x_i$ , and  $D = K \times d$ . Note that different from Bag of Word (BoW), which embeds the feature vector as follows:

$$\phi_{BoW}(x_i) = [\dots, 0, 1, 0, \dots] \in \mathbb{R}^K$$

where only  $j^{th}$  component of  $\phi_{BoW}(x_i)$  non-zero means that the nearest neighbor of feature  $x_i$  is visual word  $c_j$ ; VLAD considers the residual between a visual word and its nearest feature. Do et al. [29] show that VLAD is a simplified version of local coordinate coding [30], which tries to approximate a nonlinear classification function by a linear function.

From a set of embedded vector:  $\{\phi(x_i) \in \mathbb{R}^D \mid i = 1, \dots, M\}$ , we aggregate them by the sum pooling to obtain a single vector representing the image  $I$ :

$$\tau(I) = \sum_{i=1}^n \phi(x_i)$$

In literature, there are several other ways for this aggregation step [31], [32], but for simplification, we choose sum pooling and show that it can obtain good performance in practice.

One drawback of using local features in image retrieval is that the background features (e.g., trees, sky, road, etc) significantly outnumber features from informative objects (e.g., buildings). To alleviate this, we apply PCA projection, whitening and L2-normalization [33], which limit the impact of background features in the vector  $\tau(I)$ .

During the testing phase, we calculate the similarity between the query and database images using L2 distance and retrieve top  $R$  database images with smallest distances. Next, mean-shift algorithm is applied on  $R$  retrieved images over the translational part of their poses. We then select the largest cluster, and calculate the mean of translation and rotation [34], which is viewed as a noisy measurement  $z_t$  from the image query  $I_t$ .

### C. Updating particle weights and Resampling

For every particle, its weight is computed as the following:

$$w_t^{[i]} = p(z_t | s_t^{[i]}) \propto e^{-\frac{1}{2}(z_t - s_t^{[i]})^T \Sigma_o^{-1} (z_t - s_t^{[i]})} \quad (3)$$

where,  $\Sigma_o$  is a covariance matrix which describes the noise of the measurement obtained by observation encoder. Then, all particle weights are normalized to ensure their summation equal to 1:

$$\forall i, w_t^{[i]} = \frac{w_t^{[i]}}{\sum_{j=1}^n w_t^{[j]}} \quad (4)$$

Finally, we resample particles based on their weights by stochastic universal sampling [35]. This resampling step prevents the degeneracy problem, which can occur in the long-term localization scenario. Fig. 1 shows the filtering performed by our proposed method. At the first iteration, hypotheses are randomly generated. Hypotheses with small weights vanish if they are inconsistent with the noisy measurement from the observation encoder. Finally, the ambiguity is resolved, and the particles successfully track the vehicle. It is worth noting that in the example shown, the query and retrieved images are from different times and weather conditions.

### D. Overall algorithm

Algorithm 1 summarizes the proposed Monte Carlo-based visual localization method. The critical benefit of maintaining a set of particles is to leverage Monte Carlo principle to approximate the probability distribution  $p(s_t | u_{1:t}, z_{1:t})$ . As the number of particles  $N$  are sufficiently large, this set of particles are equivalent to a representation of probability density function. The motion model ensures the temporal smoothness of the vehicle's movement.

---

### Algorithm 1 Monte Carlo-based Visual Localization (MCVL)

---

```

1: function MCVL( $\mathcal{S}_{t-1}, s_t, I_t$ )
2:    $\bar{\mathcal{S}}_t = \mathcal{S}_t = \emptyset$ 
3:   Estimate  $z_t$  from image  $I_t$  (Section III-B)
4:   for  $i = 1$  to  $N$  do
5:     sample  $v_t^{[i]}$  and  $\psi_t^{[i]}$  as (1)
6:     compute  $s_t^{[i]}$  as (2)
7:     compute  $w_t^{[i]}$  as (3)
8:      $\bar{\mathcal{X}}_t = \bar{\mathcal{X}}_t + \langle x_t^{[i]}, w_t^{[i]} \rangle$ 
9:     normalize  $w_t^{[i]}$  as (4)
10:    for  $i = 1$  to  $n$  do
11:      draw  $i$  with probability  $w_t^{[i]}$  (Section III-C)
12:      add  $x_t^{[i]}$  to  $\mathcal{X}_t$ 
13:  return  $\mathcal{X}_t$ 

```

---

## IV. EXPERIMENTS

To validate the performance of our proposed VL method, we conduct experiments on synthetic as well as real datasets. For quantitative results, we report the mean and median translation and orientation errors. The translation error is calculated as the Euclidean distance  $\|c_{est} - c_{gt}\|_2$ . The orientation error  $|\alpha|$  is computed as the angular difference  $2\cos(|\alpha|) = \text{trace}(R_{gt}^{-1} R_{est})$  between estimated and ground truth camera rotation matrices  $R_{est}$  and  $R_{gt}$ .

### A. Implementation details

In the observation encoder, we extract SIFT feature at 4 different scales with region width of 16, 24, 32 and 40, over a grid with spacing of 2 pixels, visual vocabulary size is  $K = 128$ . The SIFT features are embedded and aggregated using VLAD to obtain a single vector of length 16384, which is then projected to a 4096 dimensional space via PCA, whitened and L2-normalized. For nearest neighbors search, we set  $R = 20$ .

Particles are initialized from Gaussian distribution with the mean from the noisy measurement in the first frame. The covariance matrices for initializing 3D location  $r_{t=1}^{[i]}$  and orientation  $\Omega_{t=1}^{[i]}$  respectively are  $\text{diag}([10, 10, 10]^T)$  and  $\text{diag}([0.001, 0.001, 1]^T)$ . The parameters for our method are set as the following:  $\Sigma_o = \text{diag}([5, 5, 5, 0.0001, 0.0001, 0.001]^T)$ ,  $\Sigma_v = \text{diag}([1, 1, 0.01]^T)$ ,  $\Sigma_\psi = \text{diag}([0.0001; 0.00001; 0.01]^T)$ ,  $\mu_v = [0.1, 0.1, 0.01]^T$ ,  $\mu_\psi = [0.001, 0.00001, 0.01]^T$ , where  $\text{diag}$  is a function that outputs a square diagonal matrix with the elements of input vector on the main diagonal. The number of particles are fixed to 1000. These are the default parameters for all experiments, unless otherwise noted.

We select PoseNet [8], MapNet [18] and MapNet+PGO [18] as our competitors. Following suggestion of [18], PoseNet uses ResNet-34 [36], adds a global average pooling layer, and parameterizes camera orientation as logarithm of a unit quaternion. MapNet+PGO employs pose graph optimization (PGO) to fuse the absolute poses (produced by MapNet) and relative odometry (from visual odometry) during the inference. PoseNet is implemented in Tensorflow, MapNet and



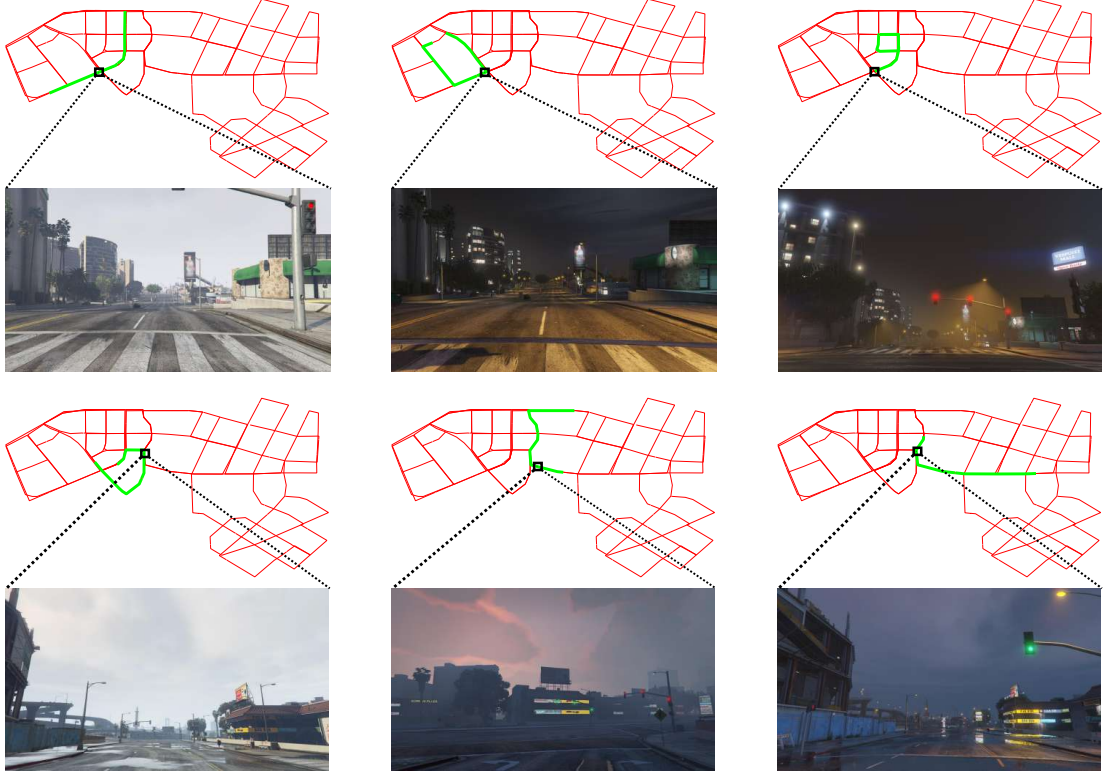


Fig. 2: Our synthetic dataset simulates image sequences collected from multiple vehicles under different environmental conditions. The red lines indicate the road network of a part of a city, while the green lines indicate a few of the trajectories taken by different data collection vehicles. Sample frames at similar locations are shown—note the differences in *pose* and environmental conditions between the images. Testing data will also be an image sequence recorded from a vehicle along the road network.

MapNet+PGO’s implementations are provided by the authors. Parameters of PoseNet and MapNet are set as the suggestion of the authors.

### B. Synthetic dataset

We construct a synthetic dataset from the computer game Grand Theft Auto V (GTA V). There are several softwares [37] [38] [39] which support researchers in collecting image data from GTA V. In this work, we adapt G2D [39] because it is able to produce an accurate ground truth 6 DoF camera pose for every image as well as provide functions to manipulate environmental conditions in the game. Note that there is an increasing recognition of the value of synthetic data towards building autonomous driving systems [40].

In particular, we simulate that there are 59 vehicles which run independently in different routes and environmental conditions. Figure 2 shows the setting of our synthetic dataset. The bird’s-eye view view of our synthetic dataset is illustrated in Figure 3, the coverage area is about 3.36km<sup>2</sup>. Also, those vehicles are simulated to run in different times of day and weather conditions. The distributions of times and weather conditions are shown in two histograms of Figure 3. The times and weathers in training sequences are uniformly distributed from 1am to 11pm, and in 7 different weather conditions

Sequences	# images	Time & Weather	Traversal distance
Test seq-1	1451	9:36 am, snowy	1393.03m
Test seq-2	360	10:41 pm, clear	359.74m
Test seq-3	1564	11:11 am, rainy	1566.93m
Test seq-4	805	6:26 pm, cloudy	806.56m
Test seq-5	1013	3:05 pm, overcast	1014.91m

TABLE I: Statistics of the testing sequences in the synthetic dataset

(snowy, foggy, clear, overcast, cloudy, sunny and rainy). Five sequences in different times and weathers are also collected for the testing phase. The statistics information and trajectory of the testing sequences are described in Table I and Figure 3 respectively. We sub-sampled on the training sequences at 2Hz and the testing sequences at 4Hz. Examples from training and testing sequences are shown in Figure 4.

In the first experiment, we compare the performance of MapNet [18] against the proposed observation encoder based on image retrieval. For our method, we use the output of the observation encoder as the final output and no filtering is done. Our purpose is to show the robustness of the observation encoder against appearance changes. For this experiment, no temporal information is used.

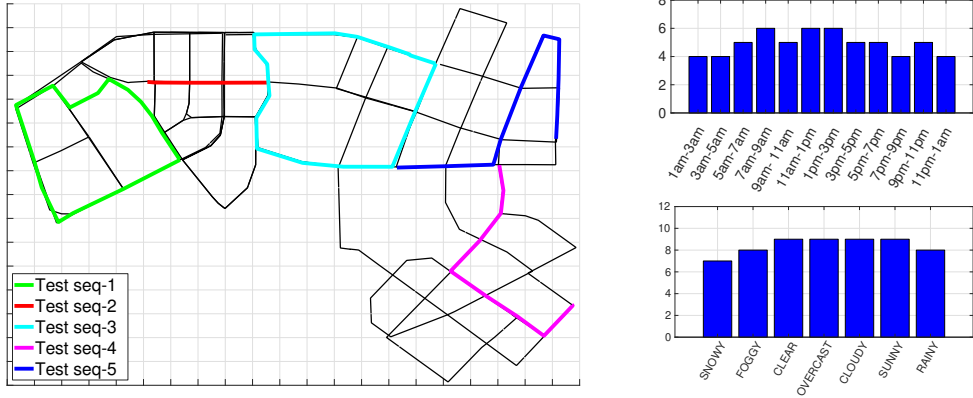


Fig. 3: Bird’s eye view and distribution of environmental conditions of our synthetic dataset. We simulate that there are 59 vehicles running in different routes, distinct time and weathers. Grid lines are every 100m. In training set, the coverage area is about  $3.36\text{km}^2$ , the time and weather conditions are uniformly distributed. The statistics of testing sequences is shown in Table I

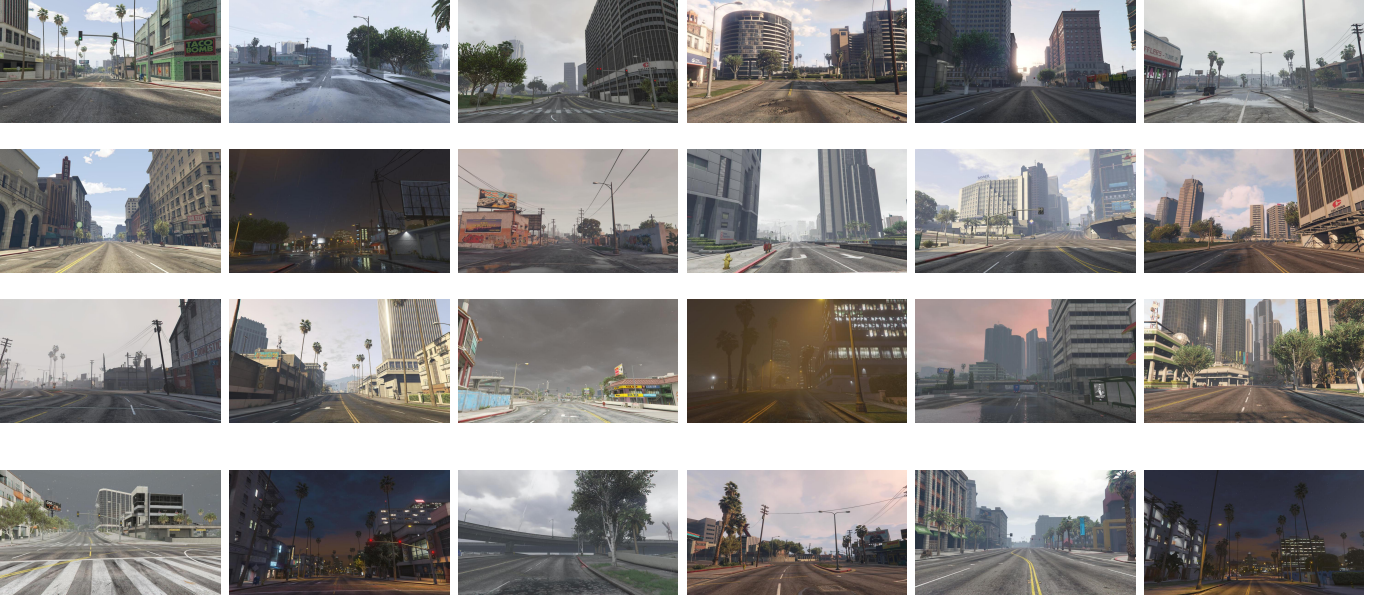


Fig. 4: Training (1st-3rd rows) and testing samples (4th row) in our synthetic dataset.

As can be seen from Table II and Figure 5, MapNet struggles to produce a reasonable result. This is because MapNet formulates the problem as an image to pose regression, whose underlying assumption of constant appearance is violated when the appearance of the environment changes. Moreover, repeated structures such as trees, sky, and road surfaces can lead to ambiguities for MapNet. In contrast, our observation encoder based on image retrieval applies state-of-the-art normalization techniques to reduce the negative impact from repetitive objects. Hence, image retrieval significantly outperforms MapNet. However, as Fig. 5 shows, using only image retrieval leads to non-smooth trajectory estimates, hinting at the need for temporal reasoning. The results for the complete system, which employs a motion model to exploit

temporal information, leads to a dramatic improvement of our method over image retrieval in terms of mean errors. The median errors are comparable between image retrieval and our method, showing that our method has a tighter error distribution.

### C. Oxford RobotCar

We compare our proposed method to state-of-the-art approaches, i.e., PoseNet [8], MapNet and MapNet+PGO [18]. In particular, PoseNet directly regresses 6-DoF camera pose from an input image. MapNet receives videos as training data, hence its loss function minimizes absolute pose per image as well as the relative pose between consecutive images, it is then followed by a fine-tuning step on unlabeled data with their visual odometry (VO). MapNet+PGO, in the inference

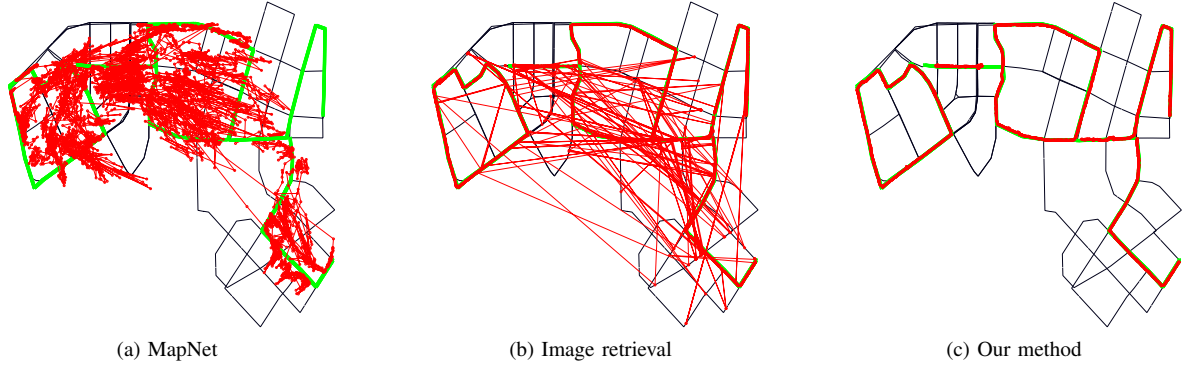


Fig. 5: Results on our synthetic dataset. Green lines represented ground truth and predicted trajectories are given in red.

	MapNet	Image retrieval	Our method		MapNet	Image retrieval	Our method
Test seq-01	71.29m, 25.14°	7.79m, 4.44°	<b>3.89m, 4.17°</b>	Test seq-01	37.45m, 4.61°	<b>2.57m, 3.31°</b>	2.63m, 3.46°
Test seq-02	42.94m, 3.25°	90.59m, 30.70°	<b>23.67m, 24.44°</b>	Test seq-02	31.06m, 0.96°	<b>4.31m, 1.38°</b>	6.12m, 3.32°
Test seq-03	138.21m, 11.88°	15.51m, 6.46°	<b>3.91m, 4.66°</b>	Test seq-03	98.34m, 4.28°	3.29m, <b>3.47°</b>	<b>3.21m, 4.03°</b>
Test seq-04	54.65m, 9.97°	6.47m, 2.31°	<b>2.97m, 2.39°</b>	Test seq-04	38.50m, 1.53°	2.73m, <b>1.17°</b>	<b>2.58m, 1.82°</b>
Test seq-05	795.41m, 11.78°	43.14m, 10.75°	<b>4.03m, 9.07°</b>	Test seq-05	807.93m, 9.71°	<b>1.78m, 7.06°</b>	1.83m, 7.29°

TABLE II: Mean (left) and median (right) rotation and translation errors on our synthetic dataset.

	Alternate route (1km)	Full route (10km)		Alternate route (1km)	Full route (10km)
PoseNet [8]	10.45m, 4.45°		PoseNet [8]	5.01m, 1.67°	
MapNet [18]	6.02m, 2.38°	28.16m, 8.50°	MapNet [18]	3.99m, 1.02°	7.08m, 1.02°
MapNet+PGO [18]	<b>5.76m, 2.27°</b>	27.75m, 8.47°	MapNet+PGO [18]	<b>3.88m, 1.01°</b>	6.99m, <b>1.01°</b>
Our method	7.15m, 3.89°	<b>5.06m, 2.48°</b>	Our method	6.74m, 1.70°	<b>4.23m, 1.40°</b>

TABLE III: Mean (left) and median (right) rotation and translation errors on the Oxford RobotCar dataset.

step, fuses the prediction of MapNet with VO by using pose graph optimization (PGO) to ensure the temporal smoothness.

We follow the configuration suggested by [18]. The split of training and testing sequences are summarized in Table IV. The training and testing sequences are recorded in different dates, which ensure a significant difference in their appearances. The experiment is conducted on the alternate and full routes with the length of 1 km and 10 km respectively. MapNet is initially trained with the training set, and then fine-tuned in an unlabeled dataset, PoseNet is trained with the training set, and our method only uses training set as the reference database.

As shown in Table III, our method outperforms PoseNet, MapNet and MapNet+PGO with a large margin in the full route, i.e., 10 km running distance. The predicted trajectory comparison between our method and MapNet and MapNet+PGO is shown in Fig. 6. Our method produces a smoother trajectory than MapNet and MapNet+PGO. In alternate route, although MapNet+PGO achieves the best quantitative result, our method outputs a smoothest trajectory, compared to other approaches. One possible reason which PGO is unable to produce a smooth prediction is there are some existing outliers from the predictions of both MapNet and VO which will probably make PGO stuck in a local minimum.

Route	Purpose	Recorded
Alternate route (1 km)	Training	26/6/2014, 9:24:58
	Training	26/6/2014, 8:53:56
	Unlabeled	14/5/2014, 13:50:20
	Unlabeled	14/5/2014, 13:46:12
Full route (10 km)	Query	23/6/2014, 15:41:25
	Training	28/11/2014, 12:07:13
	Training	02/12/2014, 15:30:08
	Unlabeled	12/12/2014, 10:45:15
	Query	09/12/2014, 13:21:02

TABLE IV: The split of training and testing sequences in Oxford RobotCar dataset.

## V. CONCLUSION

This paper presents a practical Monte Carlo-based visual localization which is shown as an efficient method to exploit the temporal smoothness of an image sequence. Inside the Monte Carlo localization scheme, we propose a novel image retrieval pipeline which represents an image as a (fixed) single vector and is robust against appearance changes. The experimental results show that our method outperforms state-of-the-art algorithms in visual localization for autonomous driving.

## REFERENCES

- [1] W. Churchill and P. Newman, "Experience-based navigation for long-term localisation," *The International Journal of Robotics Research*,



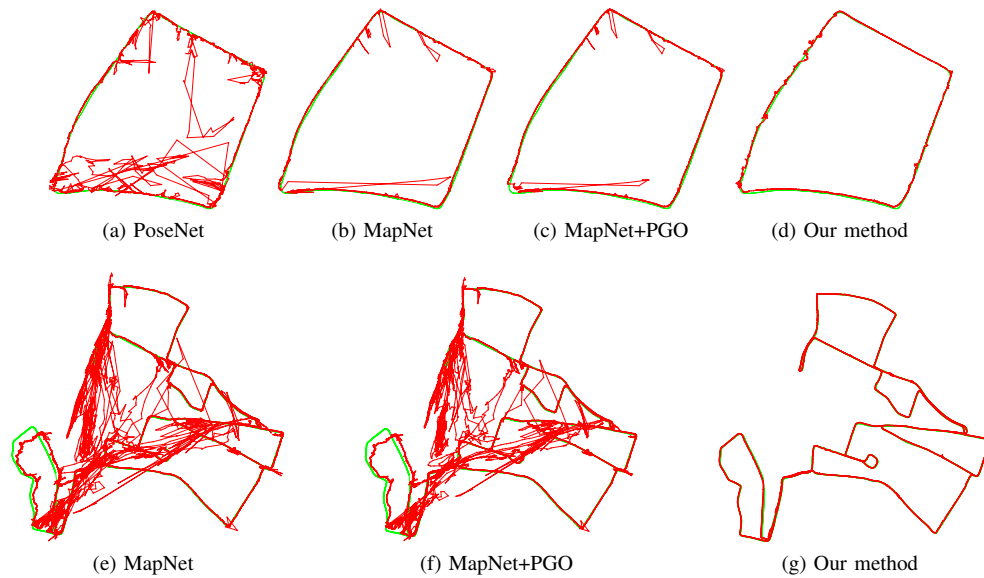


Fig. 6: Results on alternate (top) and full (bottom) routes in Oxford RobotCar dataset. The length of alternate and full routes are 1 km and 10 km respectively. The green lines are ground truth, and red lines are predicted trajectories.

- 2013.
- [2] T. Sattler, W. Maddern, C. Toft, A. Torii, L. Hammarstrand, E. Stenborg, D. Safari, M. Okutomi, M. Pollefeys, J. Sivic *et al.*, “Benchmarking 6DOF outdoor visual localization in changing conditions,” in *CVPR*, 2018.
- [3] X. Huang, X. Cheng, Q. Geng, B. Cao, D. Zhou, P. Wang, Y. Lin, and R. Yang, “The apolloscape dataset for autonomous driving,” *arXiv: 1803.06184*, 2018.
- [4] W. Maddern, G. Pascoe, C. Linegar, and P. Newman, “1 year, 1000 km: The oxford robotcar dataset,” *The International Journal of Robotics Research*, 2017.
- [5] M. J. Milford and G. F. Wyeth, “Seqslam: Visual route-based navigation for sunny summer days and stormy winter nights,” in *ICRA*, 2012.
- [6] N. Sünderhauf, P. Neubert, and P. Protzel, “Are we there yet? challenging seqslam on a 3000 km journey across all four seasons,” in *ICRA Workshop on Long-Term Autonomy*, 2013.
- [7] T. Sattler, B. Leibe, and L. Kobbelt, “Efficient & effective prioritized matching for large-scale image-based localization,” *IEEE Trans. Pattern Anal. Mach. Intell.*, 2017.
- [8] A. Kendall, M. Grimes, and R. Cipolla, “Posenet: A convolutional network for real-time 6-dof camera relocalization,” in *CVPR*, 2015.
- [9] E. Brachmann, A. Krull, S. Nowozin, J. Shotton, F. Michel, S. Gumhold, and C. Rother, “DSAC-differentiable RANSAC for camera localization,” in *CVPR*, 2017.
- [10] A. Kendall, R. Cipolla *et al.*, “Geometric loss functions for camera pose regression with deep learning,” in *CVPR*, 2017.
- [11] A. Torii, R. Arandjelovic, J. Sivic, M. Okutomi, and T. Pajdla, “24/7 place recognition by view synthesis,” in *CVPR*, 2015.
- [12] R. Arandjelovic, P. Gronat, A. Torii, T. Pajdla, and J. Sivic, “Netvlad: Cnn architecture for weakly supervised place recognition,” in *CVPR*, 2016.
- [13] V. Lepetit, F. Moreno-Noguer, and P. Fua, “EPnP: An accurate  $o(n)$  solution to the PnP problem,” *IJCV*, 2009.
- [14] N.-T. Tran, D.-K. Le Tan, A.-D. Doan, T.-T. Do, T.-A. Bui, M. Tan, and N.-M. Cheung, “On-device scalable image-based localization via prioritized cascade search and fast one-many ransac,” *TIP*, 2019.
- [15] J. L. Schönberger, M. Pollefeys, A. Geiger, and T. Sattler, “Semantic visual localization,” *CVPR*, 2018.
- [16] F. Walch, C. Hazirbas, L. Leal-Taixe, T. Sattler, S. Hilsenbeck, and D. Cremers, “Image-based localization using lstms for structured feature correlation,” in *ICCV*, 2017.
- [17] A. Kendall and R. Cipolla, “Modelling uncertainty in deep learning for camera relocalization,” *ICRA*, 2015.
- [18] S. Brahmabhatt, J. Gu, K. Kim, J. Hays, and J. Kautz, “Geometry-aware learning of maps for camera localization,” in *CVPR*, 2018.
- [19] J. Wolf, W. Burgard, and H. Burkhardt, “Robust vision-based localization by combining an image-retrieval system with monte carlo localization,” *IEEE Transactions on Robotics*, 2005.
- [20] —, “Robust vision-based localization for mobile robots using an image retrieval system based on invariant features,” in *ICRA*, 2002.
- [21] E. Menegatti, M. Zoccarato, E. Pagello, and H. Ishiguro, “Image-based monte carlo localisation with omnidirectional images,” *Robotics and Autonomous Systems*, 2004.
- [22] [https://en.wikipedia.org/wiki/Nonholonomic\\_system](https://en.wikipedia.org/wiki/Nonholonomic_system).
- [23] C. Rubino, A. Del Bue, and T.-J. Chin, “Practical motion segmentation for urban street view scenes,” in *ICRA*, 2018.
- [24] M. A. Brubaker, A. Geiger, and R. Urtasun, “Lost! leveraging the crowd for probabilistic visual self-localization,” in *CVPR*, 2013.
- [25] J. Ko and D. Fox, “Gp-bayesfilters: Bayesian filtering using gaussian process prediction and observation models,” *Autonomous Robots*, 2009.
- [26] J. L. Junkins and H. Schaub, *Analytical mechanics of space systems*. American Institute of Aeronautics and Astronautics, 2009.
- [27] R. Arandjelovic and A. Zisserman, “Three things everyone should know to improve object retrieval,” in *CVPR*, 2012.
- [28] H. Jégou, M. Douze, C. Schmid, and P. Pérez, “Aggregating local descriptors into a compact image representation,” in *CVPR*, 2010.
- [29] T.-T. Do, Q. D. Tran, and N.-M. Cheung, “Faemb: a function approximation-based embedding method for image retrieval,” in *CVPR*, 2015.
- [30] K. Yu and T. Zhang, “Improved local coordinate coding using local tangents,” in *ICML*, 2010.
- [31] H. Jégou and A. Zisserman, “Triangulation embedding and democratic aggregation for image search,” in *CVPR*, 2014.
- [32] N. Murray and F. Perronnin, “Generalized max pooling,” in *CVPR*, 2014.
- [33] H. Jégou and O. Chum, “Negative evidences and co-occurrences in image retrieval: The benefit of pca and whitening,” in *ECCV*, 2012.
- [34] F. L. Markley, Y. Cheng, J. L. Crassidis, and Y. Oshman, “Averaging quaternions,” *Journal of Guidance, Control, and Dynamics*, 2007.
- [35] D. Whitley, “A genetic algorithm tutorial,” *Statistics and computing*, 1994.
- [36] K. He, X. Zhang, S. Ren, and J. Sun, “Deep residual learning for image recognition,” in *CVPR*, 2016.
- [37] S. R. Richter, Z. Hayder, and V. Koltun, “Playing for benchmarks,” in *ICCV*, 2017.
- [38] P. Krähenbühl, “Free supervision from video games,” in *CVPR*, 2018.

- [39] A.-D. Doan, A. M. Jawaid, T.-T. Do, and T.-J. Chin, “G2D: from GTA to Data,” *arXiv preprint arXiv:1806.07381*, pp. 1–9, 2018.
- [40] J. Tremblay, A. Prakash, D. Acuna, M. Brophy, V. Jampani, C. Anil, T. To, E. Cameracci, S. Bochoon, and S. Birchfield, “Training deep networks with synthetic data: Bridging the reality gap by domain randomization,” in *CVPR Workshop on Autonomous Driving*, 2018.