

The Earth ain't Flat: Monocular Reconstruction of Vehicles on Steep and Graded Roads from a Moving Camera

Junaid Ahmed Ansari^{1*}, Sarthak Sharma^{1*}, Anshuman Majumdar¹, J. Krishna Murthy² and K. Madhava Krishna¹

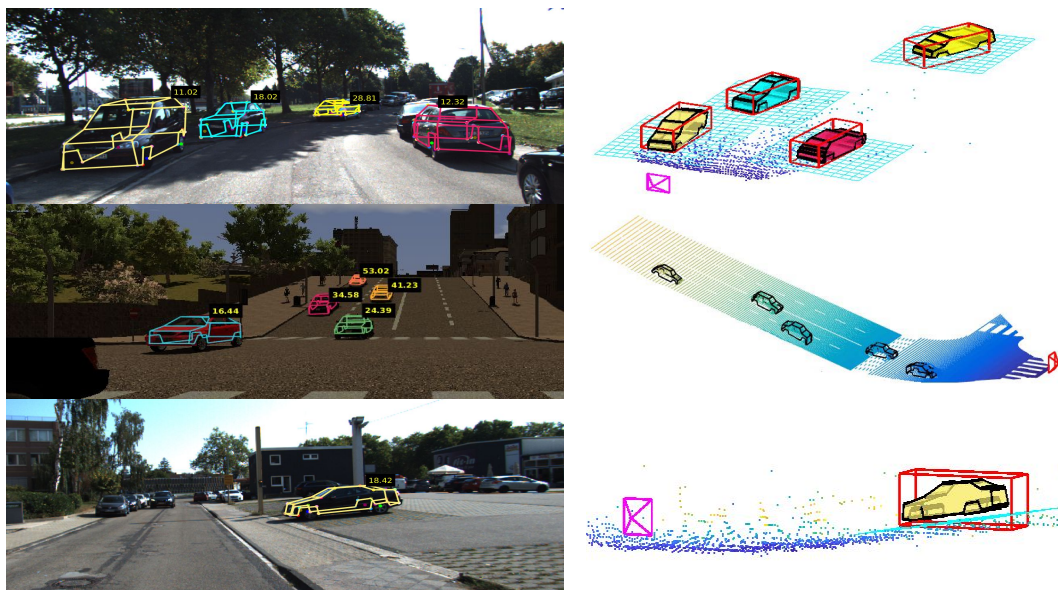


Fig. 1. Some results showcasing the efficacy of the proposed monocular object localization system. The system is capable of estimating the shape and pose (without scale-factor ambiguity) of objects located on surfaces that do not share the same plane with the moving monocular camera. The images of the scenes contain the projection of the estimated shapes (wireframes) of cars. On the top of each car, we indicate the distance of the car from the camera (in meters). To the right side of each scene, lies the visualization of the estimated wireframe and road points in 3D. For the first and third scenes, we visualize the wireframes with their respective ground truth 3D bounding boxes (shown in red) on the right, highlighting the accurate localization of the objects. In the second scene, we show the accurately estimated cars in 3D, overlaid on a dense ground truth 3D point cloud. Even the objects at over 50 meters distance on steep slopes are accurately localized.

Abstract—Accurate localization of other traffic participants is a vital task in autonomous driving systems. State-of-the-art systems employ a combination of sensing modalities such as RGB cameras and LiDARs for localizing traffic participants, but most such demonstrations have been confined to plain roads. We demonstrate, to the best of our knowledge, the first results for monocular object localization and shape estimation on surfaces that do not share the same plane with the moving monocular camera. We approximate road surfaces by local planar patches and use semantic cues from vehicles in the scene to initialize a local bundle-adjustment like procedure that simultaneously estimates the pose and shape of the vehicles, and the orientation of the local ground plane on which the vehicle stands as well. We evaluate the proposed approach on the KITTI and SYNTHIA-SF benchmarks, for a variety of road plane configurations. The proposed approach significantly improves the state-of-the-art for monocular object localization on arbitrarily-shaped roads.

*The first two authors contributed equally to this work.

¹Junaid Ahmed Ansari, Sarthak Sharma, Anshuman Majumdar, and K. Madhava Krishna are with the Robotics Research Center, KCIS, IIIT Hyderabad, India. junaid.ansari@research.iiit.ac.in

²J. Krishna Murthy is with Montreal Institute of Learning Algorithms (MILA), Université de Montreal, Canada.

I. INTRODUCTION

With the advent and subsequent commercialization of autonomous driving, there is an increased interest in monocular object localization for urban driving scenarios. While recent monocular localization methods [1], [2] achieve better localization precision when compared with stereo methods, they are confined to scenarios where the road is (very nearly) flat. This holds true for other monocular object localization systems as well [3], [4].

Reconstruction of vehicles from a monocular camera is a challenging task, owing to several factors viz. dearth of stable feature tracks on moving vehicles, self-occlusions, and it is ill-posed if the camera itself is in motion. To overcome some of these, discriminative features [5] and shape priors [2], [6] have been used to pose a bundle adjustment like scheme [2] that solves for shape and pose of a detected vehicle, assuming a prior on the shapes of all instances from a category. Using shape priors results in a richer representation of reconstructed vehicles; they are now reconstructed as 3D wireframes rather than 3D bounding boxes.

We present, to the best of our knowledge, the first results

for monocular object shape and pose estimation on surfaces that do not share the same plane with the moving monocular camera. We approximate road surfaces by local planar patches and use semantic cues from vehicles in the scene to initialize a local bundle-adjustment like procedure that simultaneously estimates the pose and shape of the vehicles, and the orientation of the local ground plane on which the vehicle stands as well. Using the proposed approach, we accurately reconstruct vehicles, predominantly using cues from only a single image. The presented method works across a variety of road geometries and demonstrate substantial improvements in terms of vehicle localization accuracy on extremely steep and non-planar roads.

To evaluate our approach, we use the popular KITTI [7] and SYNTHIA-SF [8] benchmarks. While sequences from the KITTI [7] dataset only have mild-to-moderate slopes and banks, it provides a fair comparison with other baseline methods [1], [2]. SYNTHIA-SF [8], on the other hand, has extremely steep roads and demonstrates the efficacy of the proposed approach in adapting to a wide range of road surfaces.

The remainder of the paper is organized as follows. Section II briefly discusses relevant work on monocular object localization and reconstruction in urban driving scenarios. The proposed approach is outlined in section III. In section IV, we present an evaluation of the proposed approach on popular benchmarks and discuss the results obtained thereof. Section V concludes the paper.

II. RELATED WORK

In this section, we briefly review relevant literature and contrast it with the proposed approach.

A. Shape Priors

Shape priors have been widely used in [6], [9], [2] to ease the task of object reconstruction. The underlying hypothesis is that the shape of any instance from a category can be represented as a linear combination of deformations of the mean shape for the category along certain directions, referred to as *basis vectors*. This linear subspace model was used to formulate a stochastic hill climbing problem in [6] to estimate the shape and pose of a vehicle in a single image. However, this is prohibitively slow to be used in real-time.

B. Monocular Localization in Urban Driving Scenarios

Estimating the 3D shape and pose from a single image has attracted a lot of interest in recent years, supported with the availability of datasets such as KITTI [7], ShapeNet [10] etc.

Approaches such as [1], [11] follow a 3D-2D pipeline that involves modeling the 3D shape offline and then solving for the 3D deformations in that shape using localized 2D keypoints in RGB image as evidence, thus overcoming the need to explicitly estimate the 3D keypoints. In [1], an approach to estimate the 3D shape and pose of the vehicles from a single image is presented. The 3D shape of an instance was modeled using a shape prior based on a linear

subspace model and deformation coefficients were estimated by solving an optimization problem using vehicle keypoints localized in 2D using a CNN.

In [3], [12] the authors develop a real-time monocular SfM system leveraging information from multiple image frames. However vehicles are represented as 3D bounding boxes. It was demonstrated in [2] that having a richer representation for the vehicle, such as a 3D wireframe, significantly boosts localization accuracy. Mono3D [13] trains a CNN that jointly performs object detection in 2D and in 3D space and estimates oriented bounding boxes for vehicles. Although it outperforms stereo competitors, it made the assumption of a planar road surface.

Similarly, [1], [3], [12] rely on the assumption that the plane of the vehicle to be localized is coplanar with the plane of the ego car. Most of these methods use the approach outlined in [4] to estimate the depth to a vehicle under the co-planarity assumption.

C. Monocular Road Surface Reconstruction

There is relatively little work on road surface estimation from a monocular camera. In [14], the authors propose a simple road edge prediction framework using edges and lanes detected in earlier frames. No surface level reconstruction is provided. In [15], road width and shape of the drivable area are estimated using a Conditional Random Field (CRF).

In contrast to the above approaches, the proposed approach is independent of the road plane profile and accurately localizes the vehicle independent of its coplanarity with the ego vehicle. The cost functions provided are robust, fast, and easy to implement; resulting in very accurate shape and pose estimation of the vehicle independent of the plane on which the vehicle is located. The method outperforms the current best competitor [2] by a significant margin, highlighting how the existing approaches fail when presented with non-planar road surfaces.

III. GEOMETRY AND OBJECT SHAPE COSTS

In this section, we outline our approach to reconstruct vehicles on arbitrarily oriented roads surfaces.

A. Background: Shape Priors

Along the lines of [6], [1], [2], we assume that each vehicle (in this case, a car) is represented in 3D by a wireframe consisting of K vertices (we use $K = 36$, according to the setup in [6]), each of which has a unique semantic meaning. For instance, these vertices could be locations of headlights, tail lights, wheel centers, rooftop corners, etc. that are easily identifiable across all cars. We use a set of aligned 900 CAD models of cars from the ShapeNet [10] repository and annotate each of them with K keypoint locations in 3D. We then use the render pipeline presented in [16] to synthesize a dataset comprising about 2.4 million images of rendered cars with annotated 2D keypoint locations. Over this dataset, we train a keypoint localization network based on the stacked-hourglass architecture [5]. We use this CNN, trained entirely on synthetic data, across all experiments reported in this

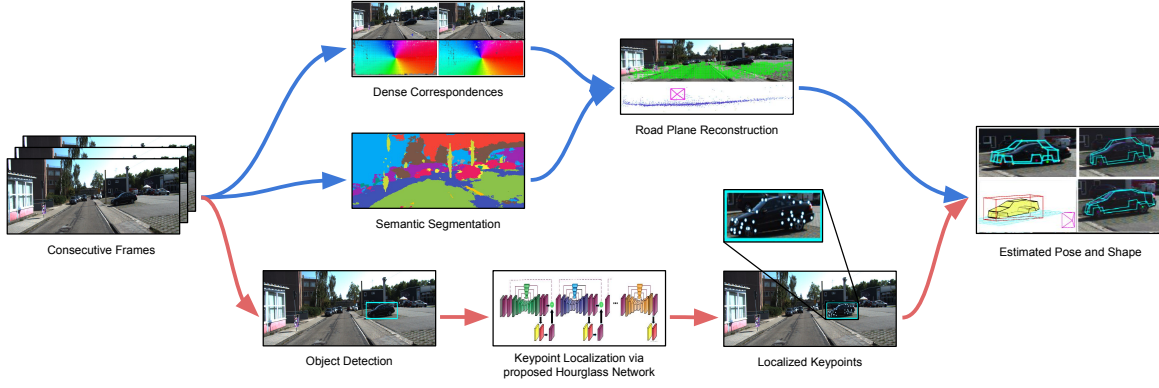


Fig. 3. Illustration of the proposed pipeline. The system takes as an input, 3 consecutive frames (in case of no lane markers). In the upper half (blue arrows), we illustrate the method for estimating the ground plane i.e. using dense correspondences over the frames and then performing bundle adjustment. In the lower half (red arrows), the detected bounding boxes in each frame are processed using the proposed keypoint localization CNN to obtain 2D locations of a discriminative set of semantic parts. The pose and shape of the object are then adjusted by incorporating the estimated ground plane information.

work. We observe that the network generalizes well to real data, consistent with the findings in [17].

Using notation from [2], we denote the mean wireframe for the vehicle category by $\bar{X} \in \mathbb{R}^{3K}$. The basis vectors are stacked into a $3K \times B$ matrix denoted V . The deformation coefficients (also referred to as the shape parameters) $\Lambda \in \mathbb{R}^B$ uniquely determine the shape of a particular instance. If we assume that the object coordinate frame has a rotation $R \in SO(3)$ and translation $t \in \mathbb{R}^3$ with respect to the camera center, any instance X can then be parameterized by the shape prior model as

$$X = \hat{R}(\bar{X} + V\Lambda) + \hat{t} \quad (1)$$

Here, $\hat{R} = \text{diag}([R, R, \dots, R]) \in \mathbb{R}^{3K \times 3K}$, and $\hat{t} = (t^T, t^T, \dots, t^T)^T \in \mathbb{R}^{3K}$. $\bar{X} = (\bar{X}_1^T, \bar{X}_2^T, \dots, \bar{X}_K^T)^T$ is an ordered collection of the 3D locations of the keypoints in the mean wireframe.

If we denote the locations of an ordered collection of 2D keypoints by $\hat{x} = (\hat{x}_1^T, \hat{x}_2^T, \dots, \hat{x}_K^T)^T \in \mathbb{R}^{2K}$, the pose (R, t) and shape (Λ) of the vehicle can be obtained by minimizing the following objective function in an alternating fashion - once for pose, and once for shape.

$$\min_{R, t, \Lambda} \mathcal{L}_r = \|\pi_K(\hat{R}(\bar{X} + V\Lambda) + \hat{t}; f_x, f_y, c_x, c_y) - \hat{x}\|_2^2 \quad (2)$$

$\pi_K()$ is a vectorized version of the perspective projection operator, which takes in K 3D points and computes their image coordinates, given the camera intrinsics $\mu = (f_x, f_y, c_x, c_y)$. Specifically, π_K is the following function.

$$\pi((X, Y, Z)^T; \mu) = \begin{pmatrix} \frac{f_x X}{Z} + c_x \\ \frac{f_y Y}{Z} + c_y \end{pmatrix}$$

$$\pi_K((X_1^T, \dots, X_K^T)^T; \mu) = (\pi(X_1; \mu)^T, \dots, \pi(X_K; \mu)^T)^T \quad (3)$$

B. System Setup

We operate on image streams captured by a front-facing monocular (RGB) camera mounted on a car. The height H above the ground at which the camera is assumed to be known a priori (this helps in resolving scale-factor ambiguity in monocular reconstruction).

We assume that, on each incoming image, an object detector [18] runs and detects vehicles in the image (as bounding boxes). We also perform a semantic segmentation of the input image using the SegNet [19] convolutional architecture. The proposed pipeline is illustrated in Fig3

C. Reconstruction of Vehicles on Slopes

To formulate a lightweight, yet robust optimization problem for reconstructing vehicles on non-planar road surfaces (i.e. roads with slopes and banks), we assume that the road is locally planar. By this, we mean that the patch of the road that lies exactly beneath a detected vehicle is assumed to be a planar patch. This assumption is corroborated by [3], where allowing each vehicle to have an adaptive local ground plane boosts localization accuracy.

Each detected vehicle v is on a planar patch parameterized by (n_g^v, d_g^v) , where n_g^v is a vector that denotes the normal to the planar patch and d_g^v denotes the distance of the planar patch from the origin of the world coordinate frame.

Resolution of Scale-Factor Ambiguity

Monocular camera setups inherently suffer from scale-factor ambiguity, i.e., any 3D length estimated from a set of images is accurate up to a positive scalar. But, for the autonomous driving applications, we require that vehicles are localized in *metric scale*, i.e., in real-world units (such as meters, for instance). We resolve scale ambiguity using one of the following two approaches.

Using Dimensions of Detected Lanes: Most roads have lane marking or zebra crossings of standard dimensions that are known to us a priori. We use the method from [20] to detect lane markings, and if we know the height

of the camera above the ground and the dimensions of the lane markings, we can retrieve the planar patch comprising the lane marking and the distance to that lane marking (in meters). Such a method estimates the local ground plane (of a lane marking near the vehicle) using information from just a single image.

Using 3-View Reconstruction and Camera Height: The above method can only be employed on roads where there are lane markings and in particular only if a lane marking is detected near a vehicle, which is not true for all scenarios we encounter. In the more general case, we can recover absolute (metric) scale by using the following 3-view reconstruction scheme. Assume we have three consecutive frames f_1, f_2, f_3 with sufficient parallax. We use DeepMatching[21] for establishing dense correspondences between frames f_1 to f_2 . Then, using a sufficient mix of road and non-road points, we estimate the egomotion between the frames using standard multi-view motion estimation techniques [22]. Using the estimated egomotion, we triangulate points *close*¹ to the car that lie on the road surface and add points from frame f_3 to the reconstruction². A local ground plane patch can then be estimated by estimating a dominant plane from the obtained point cloud using a RANSAC-like routine. Once such a plane is obtained, we can scale the reconstruction such that the median of the Y-coordinates of the estimated plane is roughly equal to the height of the camera above the ground (which is assumed to be known during initial setup).

Why

We need coplanarity to get initial pose estimate?

How does the Ground Plane help?: In scenarios where the plane of the vehicle is not same as the plane of the ego car, current methods of estimating shape and pose of the vehicle suffer due to their co-planarity assumption. We circumvent this failure by estimating the ground plane on which the vehicle is located as proposed in Fig. 4. The estimation of the ground plane parameters not only helps in correct initialization of the car, but also helps in correct localization of the vehicle constraining it to move on the plane, rather than in the line of sight of the camera to minimize the re-projection error.

Joint Optimization for Ground Plane and Vehicle Pose and Shape Estimation

Equation 2 represents the optimization problem that is solved to estimate the shape and pose of a vehicle from just a single image or from a pair of images whenever available [2]. However, this formulation assumes co-planarity of the ego car and of the object being reconstructed. We illustrate in Fig. 4 that drastic errors in localization result when the assumption does not hold.

We assume that, in the current frame, a set of vehicles \mathcal{V} have been detected by the object detection network [18]. For a particular vehicle $v \in \mathcal{V}$, we let X_i^v denote the coordinates

¹We expand the car bounding box by a factor of 1.9 to 2.0, and pick all points from the expanded bounding box that are classified as *road* by SegNet [19].

²This is typically done by propagating feature matches from frame f_2 to frame f_3 , and running a resection routine to estimate the egomotion between frame f_1 and frame f_3 , and then triangulating points from f_3 onto the initial reconstruction [22]

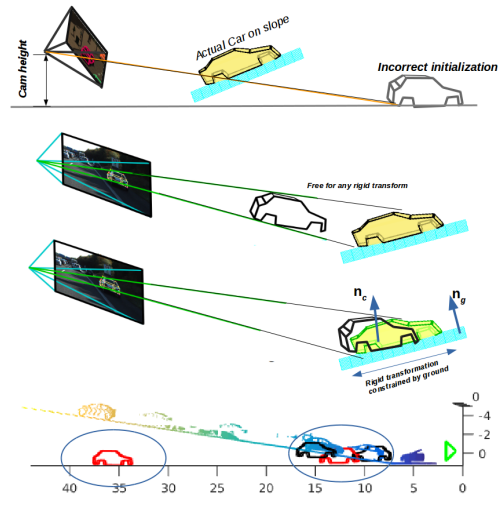


Fig. 4. From top to bottom - (i) Illustrating how co-planarity assumption results in incorrect initialization in existing approaches (ii) Relying only on minimizing the reprojection error, leaves the optimizer free to rigidly transform the mean car (iii) Joint optimization constrains the car to be on ground while minimizing the reprojection error, resulting in more accurate reconstruction and localization (n_c and n_g are car base and road plane normals respectively) (iv) Failure of co-planarity assumption for steep roads on SYNTHIA-SF [8]. Notice the incorrect initialization of the car on slopes via method of [1] shown in red. Our method is not bound by this co-planarity assumption and initializes the vehicle correctly, shown in black. We overlay the initialized wireframe on the ground truth 3D points for comparison.

of the i^{th} keypoint of the vehicle in 3D. Also, we parametrize the local ground plane beneath v by its normal vector n_g^v and the distance of the plane from the origin d_g^v . Also, we denote by n_c^v the normal of the car. The normal of the car is defined as the normal of a plane that *best*³ fits the keypoints corresponding to the wheel centers of the cars.

We now formulate a set of cost functions that relax the co-planarity assumptions in [1], [2] and estimate the vehicle’s pose and shape as well as the equation of the ground plane patch beneath it.

Ground Plane Estimation: We define a ground plane estimation loss term, which *encourages* the vehicle to be as close to the ground plane as possible. Specifically, we obtain the translation vector t_c^v to the bottom of the vehicle v^4 from the world origin (typically the camera center). This obtained quantity, in the ideal setting, represents the position vector of a point on the ground plane, the points of which are denoted as X_g^v . Formally, this term (for all vehicles in the image) can be represented as follows.

$$\mathcal{L}_g = \sum_{v \in \mathcal{V}} \|n_c^v \cdot t_c^v - d_g^v\|^2 \quad (4)$$

Normal Alignment: The normal alignment loss term stipulates that the normal of the vehicle (n_c^v) must be encouraged

³Although, in practice, all 4 wheel centers of a car are co-planar, it may still be numerically hard to determine a plane equation that satisfies all 4 points. So, we fit a plane in the least squares sense to the 4 wheel centers.

⁴We first obtain the rigid-body transform to the origin of the vehicle coordinate frame, and then concatenate to it the rigid-body transformation from the origin of the vehicle coordinate frame to the bottom of the vehicle.

to be parallel to the normal of the estimated ground plane. An initial guess for the ground plane normal is obtained as described earlier, using either lane markings, or a 3-view reconstruction. This loss can be denoted as follows. $\times(\cdot, \cdot)$ denotes the vector cross product.

$$\mathcal{L}_n = \sum_{v \in \mathcal{V}} \|\times(n_c^v, n_g^v)\|^2 \quad (5)$$

Disambiguation Prior: The above loss term has one drawback in that, it is minimized even when the estimated ground plane and vehicle normals are anti-parallel. To disambiguate such unwarranted solutions, we make use of the fact that even the steepest roads in the world have slopes less than 25 deg [23]. Whenever multiple solutions are available, we encourage the solution that's *more upright* to have a lower cost. If e_2 denotes the Y-axis of the camera coordinate system (i.e., the axis vertically pointing down), we formulate the disambiguation prior as follows (ϵ is a tiny positive constant that provides numerical stability).

$$\mathcal{L}_d = \sum_{v \in \mathcal{V}} \left\| \frac{-1}{e_2 \cdot n_c^v + \epsilon} \right\|^2 + \left\| \frac{-1}{e_2 \cdot n_g^v + \epsilon} \right\|^2 \quad (6)$$

Base Point Priors: We also use a loss term that encourages points along the base of the car (this includes keypoints on the car wheel centers, bumpers, etc) to lie as close to the estimated ground plane as possible. If X_b is a keypoint on the car base, and \mathcal{K}_b denotes the set of all keypoints that lie along the base of the car, base point priors are imposed using the following expression.

$$\mathcal{L}_b = \sum_{v \in \mathcal{V}} \sum_{X_b \in \mathcal{K}_b} \|n_c^v \cdot t_c^v - n_c^v \cdot X_b\|^2 \quad (7)$$

Global Consistency: Although we assume that each vehicle has its own planar ground patch, it is safe to assume that road planes are not susceptible to abrupt change. This is encoded into the global consistency loss term, that encourages the planar ground patch of a vehicle to be consistent with that of other vehicles around it. If \mathcal{V}^n denotes the set of all vehicles within a distance d around vehicle v (v is usually chosen to be 5 – 7 meters), the global consistency loss term is as follows.

$$\mathcal{L}_c = \sum_{v \in \mathcal{V}} \sum_{v^n \in \mathcal{V}^n} \|n_g^v - n_g^{v^n}\|^2 + \|d_g^v - d_g^{v^n}\|^2 \quad (8)$$

Dimension Regularizers: We also place priors on dimensions of vehicles that we observe, which provides a well-conditioned problem to work with and leads to better convergence rates. We use regularizers similar to ones proposed in [2], and denote the loss term by \mathcal{L}_{reg} .

Overall Optimization Problem: The overall minimization problem involving all the energy terms can be posed as follows (cf. Eq 2 4 5 6 8 7).

$$\min_{R, t, \Lambda, n_g^v, d_g^v, n_c^v} \mathcal{L}_{total} = \eta_r \mathcal{L}_r + \eta_g \mathcal{L}_g + \eta_n \mathcal{L}_n + \eta_d \mathcal{L}_d + \eta_b \mathcal{L}_b + \eta_c \mathcal{L}_c \eta_{reg} \mathcal{L}_{reg} \quad (9)$$

Here, $\eta_r, \eta_g, \eta_n, \eta_d, \eta_b, \eta_c$, and η_{reg} are weighing factors that control the relative importances of each of the loss terms. In practice, η_r, η_g, η_d , and η_b are more dominant compared to the other terms. The actual values of these weighing factors do not really matter as long as the above terms are properly weighted.

The above problem is minimized using Ceres Solver [24], a nonlinear least squares minimization framework, using a Levenberg-Marquardt optimizer with a Jacobi preconditioner.

In addition, each term is composed with a Huber loss function, to reduce the effect of outliers on the solution obtained.

Is the Huber loss important? There was no mention of it earlier.

IV. EXPERIMENTS AND RESULTS

We perform a thorough quantitative and qualitative analysis of our approach on challenging sequences from the KITTI Tracking [7] and SYNTHIA-SF [8] benchmarks. These sequences are chosen such that they capture a diverse class of road plane profiles viz. uphill, downhill, combinations of them, and even banked road planes. We compare the 3D localization error of the proposed method with the current state-of-the-art monocular competitor [2], and demonstrate significant improvements. Through a series of systematic evaluations, we demonstrate that ground plane estimation is vital for accurate localization on roads surfaces with pitch and banks. We also demonstrate that our method is independent of the road plane profile on which vehicles are to be reconstructed. In other words, unlike others (such as [1], [3], [13]) we do not make any assumptions that the ego car and the car to be reconstructed are on the same road plane.

Dataset: We use the KITTI [7] tracking benchmark to evaluate our proposed method. Sequences numbered 1, 3, 7, 8, 9, 10, 11 and 20, which contain a large number of vehicles located on roads with varying plane profiles, were used for evaluating our approach. But, KITTI [7] has only a limited number of steep slopes and banks. So, we also select about 200 vehicles located on challenging plane profiles from sequences numbered 1, 2, 4, 5 and 6 of the SYNTHIA-SF [8] dataset. We evaluate the previous best monocular competitor [2] is also evaluated on the same sequences, to ensure fair comparison.

Keypoint Network Training: The proposed network was trained on the Torch framework [25], with data comprising about 2.4 million images, generated synthetically using the modified render pipeline presented in [16]. For training and validation respectively, the generated data was split in a 75-25 ratio. The keypoint network was trained for 7 epochs on NVIDIA GTX TITAN X GPUs, spanning over about 36 hours.

A. Localization Accuracy

To evaluate localization precision, we compute the mean Absolute Translational Error (ATE) of the vehicles (in meters) of the approaches considered against the available ground truth information. We present these results in Table I, Table II and Table III. While Table I captures the overall

What is even the point of -1 within L2

TABLE I

MEAN LOCALIZATION ERROR (STANDARD DEVIATION IN PARENTHESIS) IN METERS FOR THE VEHICLES EVALUATED USING OUR APPROACH ON THE KITTI [7] TRACKING DATASET (HERE ($<x\text{ m}$) AND ($>x\text{ m}$) DENOTE THE SET OF ALL CARS WITHIN A GROUND-TRUTH DISTANCE OF x METERS AND BEYOND THE DEPTH OF x METERS RESPECTIVELY)

Approach	Overall (m)	$\leq 15\text{m}$	$\leq 30\text{m}$	$>30\text{m}$
Murthy et. al. [2]	2.61 (± 2.23)	1.59 (± 0.96)	2.52 (± 2.16)	4.30 (± 2.83)
Ours (with co-planarity assumption)	1.00 (± 0.77)	0.67 (± 0.50)	0.94 (± 0.69)	2.19 (± 1.18)
Ours (joint optimization)	0.86 (± 0.87)	0.55 (± 0.50)	0.79 (± 0.79)	2.16 (± 1.18)

TABLE II

MEAN LOCALIZATION ERROR (STANDARD DEVIATION IN PARENTHESIS) IN METERS FOR THE VEHICLES WITH CHALLENGING ROAD PROFILES EVALUATED USING OUR APPROACH ON THE KITTI [7] TRACKING DATASET

Approach	Overall (m)	$\leq 15\text{m}$	$>15\text{m}$
Murthy et. al. [2]	2.55 (± 3.16)	2.32 (± 2.21)	2.92 (± 3.38)
Ours (with co-planarity assumption)	0.95 (± 0.89)	0.92 (± 0.68)	1.00 (± 0.96)
Ours (joint optimization)	0.67 (± 0.66)	0.64 (± 0.60)	0.72 (± 0.71)

TABLE III

MEAN LOCALIZATION ERROR (STANDARD DEVIATION IN PARENTHESIS) IN METERS FOR THE VEHICLES (INCLUDING CHALLENGING ROAD PROFILE) EVALUATED USING OUR APPROACH ON THE SYNTHIA-SF [8] DATASET

Approach	Overall (m)	$\leq 15\text{m}$	$\leq 30\text{m}$	$>30\text{m}$
Murthy et. al. [2]	76.34 (± 94.03)	54.21 (± 47.93)	66.28 (± 88.74)	86.40 (± 99.32)
Ours (with co-planarity assumption)	32.03 (± 45.60)	6.3 (± 19.17)	21.76 (± 65.76)	42.31 (± 25.42)
Ours (joint optimization)	0.92 (± 0.93)	0.66 (± 0.49)	0.82 (± 0.76)	1.23 (± 1.11)

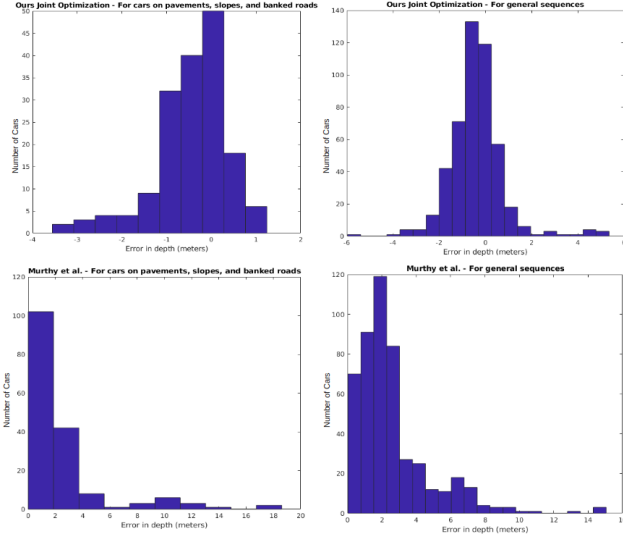


Fig. 5. Histogram showing the distribution of localization errors. Top (left to right): Plot for sequences from KITTI [7] that exhibit road slant, pitch, and banking. Plot for all evaluated sequences from the KITTI [7] benchmark. These plots show the performance of the proposed approach. The bottom plots are identical, but show the performance of the approach proposed in Murthy et. al. [2]

performance of our approach on KITTI [7] dataset, Table II presents an analysis of the performance of our approach on KITTI sequences with cars on roads with some pitch or banking angle, or parked on pavements. In Table III, we perform a thorough analysis of our approach on SYNTHIA-SF [8] which has extremely steep roads, and demonstrate the efficacy of the proposed approach in adapting to a wide variety of road plane profiles.

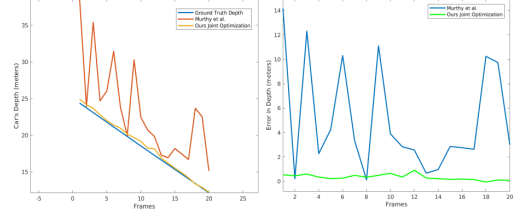


Fig. 6. Left: Predicted depth of a car on a steep slope. We compare predictions with our method with those from [2] against the ground truth. Right: Localization error for the same car when using the proposed method and when using [2].

We outperform the current best monocular localization result of [2] on the KITTI benchmark by a significant margin. It is important to note that in [2], the shape priors comprised 14 keypoints per vehicle, whereas we use a different shape prior model comprising 36 keypoints per vehicle. However, to emphasize that this improvement does not stem from more expressive shape prior used in this work, we re-implement the approach in [2] using our learnt shape priors and provide an ablation study to further drive the point home. This highlights the importance of the inclusion of ground plane in localization. As shown in Table I, we achieve a mean localization error of 0.86 meters, as compared to 2.61 meters in [2]. This is a mark improvement stemming from the inclusion of ground plane.

We also address the challenging sequences with moderate slopes on KITTI and provide our localization errors in Table II, and perform an ablation study of our approach to highlight how our the inclusion of ground plane reduces the localization error to 0.67 meters, as compared to an error of 2.55 meters given by [1]. The current state-of-the-art [1]

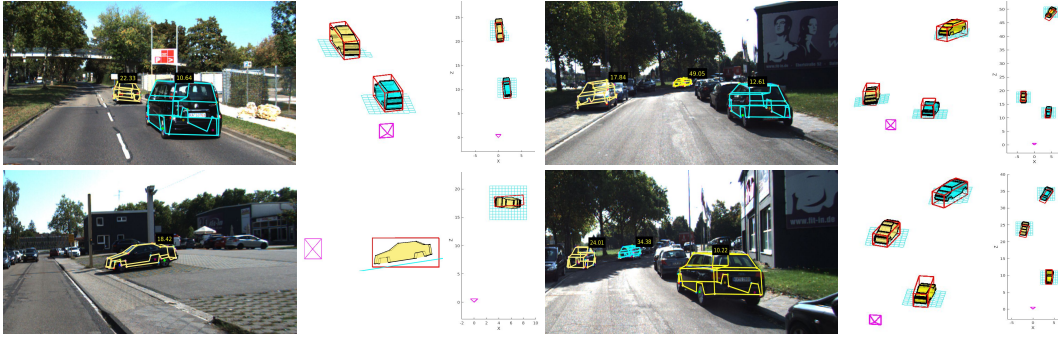


Fig. 7. Qualitative results on KITTI (with static and moving cars). The images of the scenes contain the projection of the estimated shapes (wireframes) of the cars. On the top of each car, it's depth w.r.t. the camera is displayed. Beside each scene, lies the visualization of the estimated wireframe in 3D, and the bird's eye view of the cars, along with the overlaid ground truth bounding box (in red).

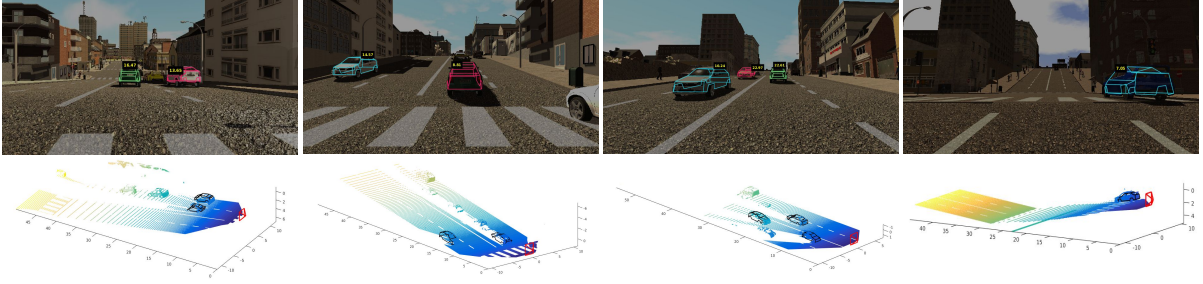


Fig. 8. Qualitative results on SYNTHIA-SF (with static and moving cars). The images of the scenes contain the projection of estimated shapes (wireframes) of the cars. On the top of each car, depth w.r.t. the camera is displayed. Below each scene image, lies the visualization of the estimated wireframe in 3D, overlaid on a dense 3D point cloud. The proposed method is able to generalize well on cars present on different plane profiles.

relies on the assumption that the plane of the vehicle and the ego car is same, i.e they are co-planar. We circumvent this assumption leading to a highly accurate localization of the vehicle irrespective of the fact that it is co-planar with the ego car or not. For vehicles that are close to the car, we achieve a high degree of precision (mean error of about 0.67 meters, with a low standard deviation as well).

To further evaluate our approach, we test it on the extremely challenging SYNTHIA-SF [8] dataset which has steep road surfaces, having various non-planar profiles. [1] fails completely in the task of accurate localization of objects in such scenarios, due to the assumption that the plane of the vehicle and that of the ego car is same, and fails to recover the correct shape and pose. Moreover, the method given by [2] fails drastically in non-planar surfaces, giving a mean localization error of 76.34 meters, amplified by the non-generalizable nature of the 14 keypoint network which leads to inaccurate keypoint localizations. Our system achieves a mean localization error of 0.92 meters, the results of which are shown in Table III. The proposed method generalizes well to different plane profiles and performs significantly better as compared to the approach of [1], which assumes coplanarity of the vehicles and ego car and hence fails in such challenging road profiles. Fig. 5 shows the error distribution of our approach (top two) and for [1] (bottom two).

B. Keypoint Localization

To evaluate the accuracy of our 2D keypoint localization network, we use the standard PCK (Percentage of Correct Keypoints) and APK (Average Precision of Keypoints) metrics, used in [26], [5] and [27]. A very tight threshold of 2 pixels is used in our experiments and analysis, for the determination of the correctness of our keypoint estimate. Our trained keypoint model achieved a PCK measure of **96.89%** at $\alpha = 0.1$ APK, on the aforementioned validation set. The network was deployed on KITTI and SYNTHIA-SF datasets.

C. Qualitative Results

We showcase the qualitative results of our approach on challenging KITTI and SYNTHIA-SF scenes with moderate to high slopes. For KITTI, in Fig. 7, we overlay the final estimate of the car in 3D along with the ground truth 3D bounding box to show how our approach estimates the vehicle shape and pose accurately. For SYNTHIA-SF, in Fig. 8, we overlay the estimate of the car after shape and pose adjustment on the ground truth scene points to highlight the accurate shape and pose estimation of the car.

D. Summary of Results

The cornerstone of this effort was to highlight that the presence of non-planar road profiles leads to an unsuccessful pose estimation of cars in urban scenarios by the current state-of-the-art approach, due to the fact that it relies on the

co-planarity of the ego car and the vehicle. Our proposed approach is independent of the plane profile on which the car is located. We improve by a large margin through the inclusion of the ground plane in KITTI sequences, which have moderate slopes. We report these results in Table I and in Table II. The importance of the proposed approach is highlighted in Table II, where we achieve a performance boost of about 4 times in scenes with moderate slopes. For an overall comparison on KITTI, we evaluate our approach on scenes with different planar and non-planar road surfaces and show an improvement of about 3 times. We further present the performance of our approach on SYNTHIA-SF [8] which has extremely steep scenes, resulting in a catastrophic failure of the current state-of-the-art monocular shape and pose estimation [1]. Our performance is significantly improved in such scenes, irrespective of the road profiles, the results of which are reported in Table III. We also perform an ablation study, reported in Table I, Table II and Table III, to highlight the importance of our ground plane estimation policy, and show that it provides a significant performance boost over just the utilization of a well constrained 36 keypoint system.

V. CONCLUSIONS

In this work, we presented an approach for accurate 3D localization and shape estimation of vehicles on steep road surfaces. Most current monocular localization systems assume co-planarity of the vehicle to be localized and the ego car, for accurate localization. However, since the assumption does not always hold in the real world, we propose the incorporation of ground plane information (and joint estimation of that information). We show that, this works well in practice, as evident by significant improvements over the state-of-the-art monocular localization methods and thus make a strong case for exploiting ground plane information. Future work could work towards building denser models of roads, and focus on heavy traffic situations - where not much of the road surface is visible.

REFERENCES

- [1] J. K. Murthy, G. S. Krishna, F. Chhaya, and K. M. Krishna, "Reconstructing vehicles from a single image: Shape priors for road scene understanding," in *Proceedings of the IEEE Conference on Robotics and Automation*, 2017.
- [2] J. K. Murthy, S. Sharma, and K. M. Krishna, "Shape priors for real-time monocular object localization in dynamic environments," in *Proceedings of the IEEE Conference on Intelligent Robots and Systems (In Press)*, 2017.
- [3] S. Song and M. Chandraker, "Joint sfm and detection cues for monocular 3d localization in road scenes," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2015.
- [4] G. P. Stein, O. Mano, and A. Shashua, "Vision-based acc with a single camera: Bounds on range and range rate accuracy," in *Intelligent Vehicles Symposium*. IEEE, 2003.
- [5] A. Newell, K. Yang, and J. Deng, "Stacked hourglass networks for human pose estimation," in *European Conference on Computer Vision*. Springer, 2016.
- [6] M. Z. Zia, M. Stark, and K. Schindler, "Towards scene understanding with detailed 3d object representations," *International Journal of Computer Vision*, 2015.
- [7] A. Geiger, P. Lenz, and R. Urtasun, "Are we ready for autonomous driving? the kitti vision benchmark suite," in *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2012.
- [8] D. Hernandez-Juarez, L. Schneider, A. Espinosa, D. Vazquez, A. M. Lopez, U. Franke, M. Pollefeys, and J. C. Moure, "Slanted stixels: Representing san francisco's steepest streets," in *British Machine Vision Conference (BMVC)*, 2017.
- [9] S. Tulsiani, A. Kar, J. Carreira, and J. Malik, "Learning category-specific deformable 3d models for object reconstruction," *IEEE transactions on pattern analysis and machine intelligence*, 2016.
- [10] A. X. Chang, T. Funkhouser, L. Guibas, P. Hanrahan, Q. Huang, Z. Li, S. Savarese, M. Savva, S. Song, H. Su, et al., "Shapenet: An information-rich 3d model repository," *arXiv preprint arXiv:1512.03012*, 2015.
- [11] X. Zhou, M. Zhu, S. Leonardos, K. G. Derpanis, and K. Daniilidis, "Sparseness meets deepness: 3d human pose estimation from monocular video," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 4966–4975.
- [12] S. Song and M. Chandraker, "Robust scale estimation in real-time monocular sfm for autonomous driving," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2014, pp. 1566–1573.
- [13] X. Chen, K. Kundu, Z. Zhang, H. Ma, S. Fidler, and R. Urtasun, "Monocular 3d object detection for autonomous driving," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2016, pp. 2147–2156.
- [14] F. Chausse, R. Aufrere, and R. Chapuis, "Recovering the 3d shape of a road by on-board monocular vision," in *Proceedings 15th International Conference on Pattern Recognition. ICPR-2000*, 2000.
- [15] J. Fritsch, T. Kühnl, and F. Kummert, "Monocular road terrain detection by combining visual and spatial information," *IEEE Transactions on Intelligent Transportation Systems*, 2014.
- [16] J. K. M. K. M. K. Parv Parkhiya, Rishabh Khawad and B. Bhowmick, "Constructing category-specific models for monocular object slam," in *Proceedings of the IEEE Conference on Robotics and Automation (in press)*, 2018.
- [17] H. Su, C. R. Qi, Y. Li, and L. J. Guibas, "Render for cnn: Viewpoint estimation in images using cnns trained with rendered 3d model views," in *The IEEE International Conference on Computer Vision (ICCV)*, December 2015.
- [18] J. X. J. W. J. QiongYan and Y.-W. LiXu, "Accurate single stage detector using recurrent rolling convolution," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017.
- [19] V. Badrinarayanan, A. Kendall, and R. Cipolla, "Segnet: A deep convolutional encoder-decoder architecture for image segmentation," *IEEE transactions on pattern analysis and machine intelligence*, vol. 39, no. 12, pp. 2481–2495, 2017.
- [20] R. K. Satzoda and M. M. Trivedi, "Vision-based lane analysis: Exploration of issues and approaches for embedded realization," in *Computer Vision and Pattern Recognition Workshops (CVPRW)*, 2013 *IEEE Conference on*. IEEE, 2013, pp. 604–609.
- [21] P. Weinzaepfel, J. Revaud, Z. Harchaoui, and C. Schmid, "DeepFlow: Large displacement optical flow with deep matching," in *IEEE International Conference on Computer Vision (ICCV)*, Sydney, Australia, Dec. 2013. [Online]. Available: <http://hal.inria.fr/hal-00873592>
- [22] R. I. Hartley and A. Zisserman, *Multiple View Geometry in Computer Vision*, 2nd ed. Cambridge University Press, ISBN: 0521540518, 2004.
- [23] "Kiwi climb: Hoofing up the world's steepest street," <http://edition.cnn.com/travel/article/worlds-steepest-street-residents/index.html>.
- [24] S. Agarwal, K. Mierle, and Others, "Ceres solver," <http://ceres-solver.org>.
- [25] R. Collobert, K. Kavukcuoglu, and C. Farabet, "Torch7: A matlab-like environment for machine learning," in *BigLearn, NIPS Workshop*, 2011.
- [26] S. Tulsiani and J. Malik, "Viewpoints and keypoints," in *2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE, 2015.
- [27] Y. Yang and D. Ramanan, "Articulated pose estimation with flexible mixtures-of-parts," in *2011 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE, 2011.