

Unsupervised Monocular Depth Learning with Integrated Intrinsic and Spatio-Temporal Constraints

Kenny Chen¹, Alexandra Pogue², Brett T. Lopez³, Ali-akbar Agha-mohammadi³, and Ankur Mehta¹

Abstract—Monocular depth inference has gained tremendous attention from researchers in recent years and remains as a promising replacement for expensive time-of-flight sensors, but issues with scale acquisition and implementation overhead still plague these systems. To this end, this work presents an unsupervised learning framework that is able to predict at-scale depth maps and egomotion, in addition to camera intrinsics, from a sequence of monocular images via a single network. Our method incorporates both spatial and temporal geometric constraints to resolve depth and pose scale factors, which are enforced within the supervisory reconstruction loss functions at training time. Only unlabeled stereo sequences are required for training the weights of our single-network architecture, which reduces overall implementation overhead as compared to previous methods. Our results demonstrate strong performance when compared to the current state-of-the-art on multiple sequences of the KITTI driving dataset and can provide faster training times with its reduced network complexity.

I. INTRODUCTION

Modern robotic agents take advantage of accurate, real-time range measurements to build a spatial understanding of their surrounding environments for collision avoidance, state estimation, and other navigational tasks. Such measurements are commonly retrieved via active sensors (e.g., LiDAR) which resolve distance by measuring the time-of-flight of a reflected light signal; however, these sensors are often costly [1], difficult to calibrate and maintain [2], [3], and can be unwieldy for platforms with a weight budget [4]. *Passive* sensors, on the other hand, have seen a tremendous surge of interest in recent literature to predict scene depth from input imagery using multi-view stereo [5]–[7], structure-from-motion [8]–[11], or more recently, purely monocular systems [12]–[17], due to their smaller form factor and increasing potential to rival the performance of explicit active sensors with the advent of machine learning.

In particular, monocular depth inference is attractive since RGB cameras are ubiquitous in modern times and requires the least number of sensors, but this setup suffers from a fundamental issue of scale acquisition. More specifically, in a purely monocular system, depth can only be estimated up to an ambiguous scale and requires additional geometric information to resolve the units of the depth map. Such

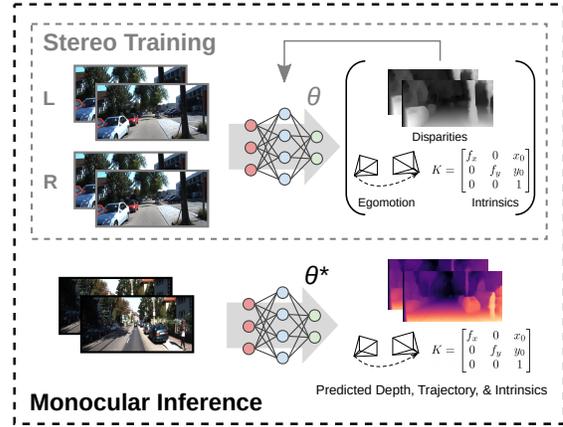


Fig. 1. **System Overview.** Our system regresses depth, pose and camera intrinsics from a sequence of monocular images. During training, we use two pairs of unlabeled stereo images and consider losses in both spatial and temporal directions for our network weights. During inference, only monocular images are required as input, and our system outputs accurately scaled depth maps and egomotion in addition to the camera’s intrinsics.

cameras typically capture frames by projecting 3D scene information onto a 2D image plane, and abstracting higher dimensional depth information from a lower dimension is a fundamentally ill-posed problem. To resolve the scale factors of these depth maps, a variety of learning-based approaches have been proposed with differing techniques to constrain the problem geometrically [13], [18]–[26]. Temporal constraints, for example, are commonly employed [12], [25], [27], [28] and is defined as the geometric constraint between two consecutive monocular frames, aiming to minimize the photometric consistency loss after warping one frame to the next. Spatial constraints [13], [19], [29], on the other hand, extract scene geometry not through a forward-backward reconstruction loss (i.e., temporally) but rather in left-right pairs of stereo images with a predefined baseline. Most works choose to design their systems around either one or the other, and while a few systems have integrated both constraints before in a multi-network framework [30]–[32], none have taken advantage of both spatial and temporal constraints in a single network to resolve these scale factors.

To this end, we propose an unsupervised, single-network monocular depth inference approach that considers both spatial and temporal geometric constraints to resolve the scale of a predicted depth map. These “spatio-temporal” constraints are enforced within the reconstruction loss functions of our network during training (Fig. 1), which aim to minimize the photometric difference between a warped frame and the actual next frame (forward-backward) while simultaneously maximizing the disparity consistency between a pair of stereo

¹Kenny Chen and Ankur Mehta are with the Department of Electrical and Computer Engineering at the University of California Los Angeles, Los Angeles, CA 90095, USA. {kennyjchen, mehtank}@ucla.edu

²Alexandra Pogue is with the Department of Mechanical and Aerospace Engineering at the University of California Los Angeles, Los Angeles, CA 90095, USA. anpogue@ucla.edu

³Brett T. Lopez and Ali-akbar Agha-mohammadi are with the NASA Jet Propulsion Laboratory, California Institute of Technology, Pasadena, CA 91109, USA. {brett.t.lopez, aliagha}@jpl.nasa.gov

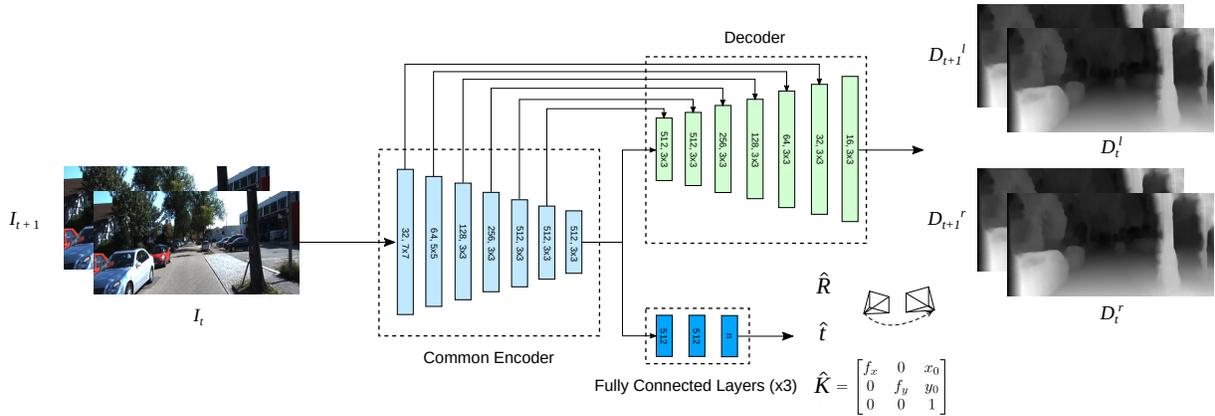


Fig. 2. **Architecture Overview.** Our system uses a common convolutional-based encoder between the different outputs, which compresses the input images into a latent space representation. This representation is then sent through either a trained decoder to retrieve left-right stereo image disparities, or through different groups of fully connected layers to estimate egomotion ($n = 3$) or camera intrinsics ($n = 4$). In the common encoder, each block uses a series of two convolutional layers, the first with stride 2 and the second stride 1 (zero padding), and with input dimensions and kernel sizes as specified. The transposed convolutional blocks in the decoder are similarly structured, with pooling indices received from the corresponding encoder’s feature maps.

frames (left-right). Unlike previous approaches, we consider camera intrinsics as an additional unknown parameter to be inferred and demonstrate accurate inference of both depth and camera parameters through a sequence of purely monocular frames; this is all performed in a single end-to-end network to minimize implementation overhead.

Our main contributions are as follows: (1) we propose an unsupervised, single-network architecture for monocular depth inference which takes advantage of the geometric constraints found in both spatial and temporal directions; (2) a novel loss function that integrates unknown camera intrinsics directly into the single-network training procedure; and (3) extensive performance and run-time analyses of our proposed architecture to verify our methods. These efforts were in support of NASA’s Jet Propulsion Laboratory’s Networked Belief-aware Perceptual Autonomy (NeBula) framework [33] as part of Team CoSTAR in the DARPA SubT Challenge.

Related Work

Depth estimation using monocular images and deep learning began with supervised methods over large datasets and ground truth labeling [22], [34], [35]. Although these methods produced accurate results, acquiring ground truth data for supervised training requires expensive 3D sensors, multiple scene views, and inertial feedback to obtain even sparse depth maps [20]. Later work sought to address a lack of available high-quality labeled data by posing monocular depth estimation as a stereo image correspondence problem, where the second image in a binocular pair served as a supervisory signal [18]–[20]. This approach trained a convolutional neural network (CNN) to learn epipolar geometry constraints by generating disparity images subject to a stereo image reconstruction loss. Once trained, networks were able to infer depth using only a single monocular color image as input. While this work achieved results comparable to supervised methods in some cases, occlusion and texture-copy artifacts that arose with stereo supervision motivated learning approaches using a temporal sequence of images

as an alternative [13], [19]. CNNs trained using monocular video regressed depth using the camera egomotion to warp a source image to its temporally adjacent target. To address the additional problem of camera pose, [12]–[17], [36], [37] trained a separate pose network.

The learning of visual odometry (VO) and depth maps has useful application in visual simultaneous localization and mapping (SLAM). Visual SLAM leverages 3D vision to navigate an unknown area by determining camera pose relative to a constructed global map of an environment. To build and localize within a map, VO in a SLAM pipeline must solve at metric scale. Geometric approaches to monocular SLAM using first principled solutions, such as structure from motion (SfM) [38], resolved scaling issues using external information [26], [39]. Building on such methods, work in data-driven monocular VO obtained scale using sources such as GPS sensor fusion [40] or training supervision [36], [39]. Unsupervised approaches using a camera alone remain attractive however, due to the reduction of manual effort associated with fewer sensors. Promising research in this area combined visual constraints (e.g., monocular depth [12]–[17], stereo depth [26], [30]–[32], or optical flow [17], [41]) to achieve scale consistent outcomes.

Network architecture for visual odometry and dense depth map estimation separate depth and pose networks into two CNNs, one with convolutional and fully connected layers and the other an encoder-decoder structure [42], respectively. In the case where only monocular images are used in training, the self-supervision inherent in estimation is less constrained, having only pose generated from temporal constraints to determine depth, and vice versa [41]. The work of [15], [43] for example, suffered from scaling ambiguity issues [30]. Training using binocular video, on the other hand, made use of independent constraints from spatial and temporal image pairs that offered an enriched set of sampled images for network training. This “spatio-temporal” approach allowed for the regression of depth from spatial cues generated by epipolar constraints, which were then passed to the

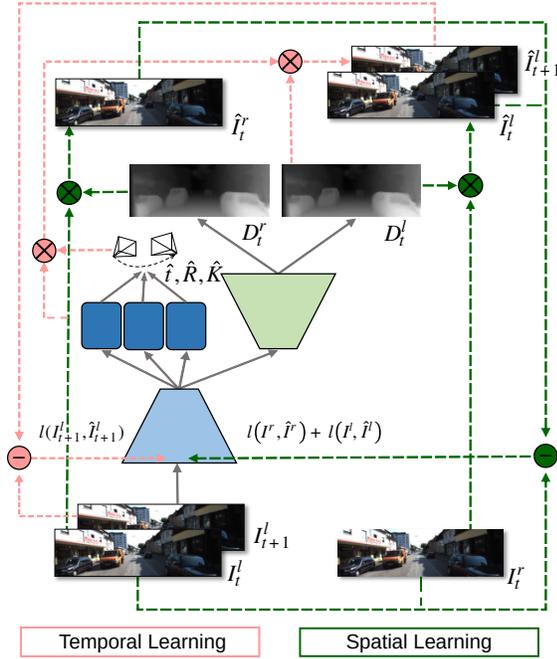


Fig. 3. **Training Diagram.** Our single-network system runs a timed sequence of left images through the common encoder (light blue trapezoid) to generate outputs that are fed to the fully connected (FC) layers (blue rectangles) and the decoder (green trapezoid). Outputs from the FC layers and the decoder are the camera pose and intrinsics, and disparity maps, respectively. The disparities are used to find left-right reprojection images (green dashed lines), while the disparities, camera pose and intrinsics determine the temporal reconstructions (pink dashed lines). All input and output images are framed in black for clarity.

pose network to independently estimate VO using temporal constraints [13], [30]–[32], [41].

In this work, we propose a spatio-temporal network inspired by [30]–[32] that uses an effective combination of losses to simultaneously regress depth and egomotion in a single network. Additionally, to provide freedom from manual calibration, the network is also capable of learning camera intrinsics which can be useful when a video source is unknown [12], [25]. To predict depth, we use a photoconsistency loss between stereo image pairs, a left-right consistency loss between image disparity maps [19], and a disparity map smoothing function [44]. To estimate egomotion and camera intrinsics, we leverage a unique loss function that accounts for the photometric difference between temporally adjacent images. By combining these losses, we show that we can obtain scaled visual odometry information and accurate camera parameters. Furthermore, by observing the similarities between the architecture of the depth network encoder and the pose network’s convolutional layers, we can effectively eliminate architecture redundancy by merging them via a common encoder (Fig. 2) and can provide faster training times without significant loss in performance.

II. METHODS

A. Notation

A color image, I , is composed of pixels with coordinates $p_{ij} \in \mathbb{R}^2$, where $I_{ij} = I(p_{ij})$. In temporal training, we

denote images at time t as I^t , and images temporally adjacent as source frame $I^{t'}$. A pixel at time t' is transformed to its corresponding pixel at time t using homogeneous transformation matrix $T_{t' \rightarrow t} \in \mathbb{SE}(3)$ and camera intrinsics matrix $K \in \mathbb{R}^{3 \times 3}$, where pixels in homogeneous coordinates, $\tilde{p} = (p, 1)^T$, are denoted p for simplicity.

Rectified stereo image pairs are given by I^r, I^l , where the superscripts for time have been dropped for convenience, and superscripts l, r correspond to the left and right images respectively. D^l represents the disparity map that warps I^r to the corresponding I^l , and we define per pixel disparity as $d_{ij}^l = D^l(p_{ij})$. Thus $I_{ij}^l = I_{i+d^l, j}^r$, and $d_{i+d^l, j}^r = D^r(p_{i+d^l, j})$ is the disparity that does the reverse operation. Depth per pixel z is then determined by the relation, $z = Bf_x/d$, where f_x is the x -component focal length and B is the horizontal baseline between stereo cameras.

B. Preliminaries

We can obtain the projected pixel coordinates and depth map using equation,

$$z^t p^t = K R_{t' \rightarrow t} K^{-1} z^{t'} p^{t'} + K t_{t' \rightarrow t}, \quad (1)$$

where the intrinsics matrix, K , is written explicitly as:

$$K = \begin{bmatrix} F & X_0 \\ 0 & 1 \end{bmatrix}, \quad F = \text{diag}(f_x, f_y), \quad X_0 = [x_0, y_0]^T, \quad (2)$$

and R and t are the rotation matrix and translation vector arguments of transformation matrix T [12]. Note that in this work we assume no lens distortion and a zero skew coefficient in the camera, and that stereo cameras have equal intrinsic parameters. Equation (1) constitutes the temporal reconstruction loss at training used to determine the camera egomotion, R and t , and the camera intrinsics K in a single network.

C. Overall Optimization Objective

Our loss function is made up of a novel temporal reconstruction term and four spatial reconstruction terms [19], [26]. Error regression for the following losses allows the network to correctly predict a target image temporally and spatially during training in order to infer depth, pose, and camera intrinsics from a monocular image sequence at test time. The temporal reconstruction term of the loss function is implicitly defined where $l_{te}(I^{l,t}, I^{l,t'}) \rightarrow \hat{I}^{l,t}$, and the spatial reconstruction terms are composed of a photoconsistency loss, l_p , a left-right consistency loss l_{lr} , and a disparity smoothness loss l_r ,

$$\begin{aligned} l(f(I^l; \theta), I^l, I^r) &= \lambda_p (l_p(I^l, \hat{I}^l) + l_p(I^r, \hat{I}^r)) \\ &+ \lambda_{te} l_p(I^{l,t}, \hat{I}^{l,t}) + \lambda_{lr} l_{lr}(D^l, D^r) \\ &+ \lambda_r (l_r(D^l, I^l) + l_r(D^r, I^r)). \end{aligned} \quad (3)$$

The argument I in the loss function is the original image and \hat{I} is the reprojected image, and individual losses are weighted by λ labeled with corresponding subscripts.

1) *Spatio-Temporal Reconstruction Loss*: The *photoconsistency loss* compares image appearance using the structural similarity index measure (SSIM) and an absolute error between generated and sampled images [19], [32], [45]:

$$l_p(I, \hat{I}) = \frac{1}{N} \sum_{i,j} \alpha \frac{1 - SSIM(I_{ij}, \hat{I}_{ij})}{2} + (1 - \alpha) |I_{ij} - \hat{I}_{ij}|. \quad (4)$$

The loss is composed of three terms in total (two spatial losses and a temporal loss). N in this equation is the number of image pixels and the weight α is set to 0.85.

For reprojected images, we assume equal camera intrinsics produce right and left stereo images. The focal length \hat{f}_x from intrinsics matrix \hat{K} in (2) is co-predicted via the learned disparity and penalized using spatial reconstruction losses. For stereo image inputs, predicted disparity maps are used to generate the left view from a right image, and vice versa. Depth values calculated from the disparity maps are then input to the temporal reconstruction loss to generate the left target image from temporally adjacent source images, i.e. to generate the temporal image arguments for (4), we put (1) in the form where for pixels $P = \{p_i, i = 1 \dots N\}$,

$$\sum_{i,j} |I_{ij}^{l,t} - \hat{I}_{ij}^{l,t}| \rightarrow \sum_{p \in P} \left| z^t p^t - \left[\hat{K} \hat{R}_{t' \rightarrow t} \hat{K}^{-1} \frac{b \hat{f}_x}{d'} p^{t'} + \hat{K} \hat{t}_{t' \rightarrow t} \right] \right|, \quad (5)$$

is the absolute error between the left image and the reprojected image, and the structure similarity measure is generated by the same mappings between I_{ij} and pixel p_i .

We distinguish this loss function from previous works in the following ways: depth estimation in the above temporal relation is derived using spatial losses from within the same network, and the temporal loss infers both egomotion and camera intrinsics. This goes beyond work that used solely temporal constraints and a separate depth network [12], [25] or spatio-temporal work that used predetermined intrinsics and a separate depth network [13], [30]–[32], [41].

2) *Spatial Reconstruction Loss*: The *left-right disparity consistency loss* is used to obtain consistency between disparity maps [19]. During training, the network predicts disparity maps D^l and D^r using only left image sequences as input and then penalizes the difference between the left-view disparity map and the warped right view, as well as the right-view and the warped left view,

$$l_{lr}(D^l, D^r) = \frac{1}{N} \sum_{i,j} |d_{ij}^l - d_{i+d^l,j}^r| + |d_{ij}^r - d_{i+d^r,j}^l|. \quad (6)$$

The *disparity smoothness loss* penalizes depth discontinuities that occur at image gradients ∂I [44]. To obtain locally smooth disparities, an exponential weighting function is used on disparity gradients ∂d :

$$l_r(D, I) = \frac{1}{N} \sum_{i,j} |\partial_x d_{ij}| e^{-|\partial_x I_{ij}|} + |\partial_y d_{ij}| e^{-|\partial_y I_{ij}|}. \quad (7)$$

D. Learning the Camera Intrinsics

For predicted parameters \hat{K} , \hat{R} , \hat{t} in (1), penalizing differences via training loss ensures $\hat{K}\hat{t}$ and $\hat{K}^{-1}\hat{R}\hat{K}$ converge to the correct values. To determine parameters individually, the translational relation fails because it is under-determined since there exists incorrect values of \hat{K} and \hat{t} such that $\hat{K}\hat{t} = Kt$. The rotational relationship, $\hat{K}\hat{R}\hat{K}^{-1} = KRK^{-1}$, however, does uniquely determine \hat{K} , \hat{R} such that they are equal to K , R , and therefore provides sufficient supervisory signal to estimate these values accurately.

Proof: From the above relation we obtain $\hat{R} = \hat{K}^{-1}KRK^{-1}\hat{K}$, and we constrain \hat{R} to be $SO(3)$, i.e. $\hat{R}^T = \hat{R}^{-1}$ and $\det(\hat{R}) = 1$. Substituting \hat{R} into the relationship $\hat{R}\hat{R}^T = I$, we find that $AR = RA$ where $A = K^{-1}\hat{K}\hat{K}^TK^{-T}$. The value $\det(\hat{K}^{-1}KRK^{-1}\hat{K})$ is equal to 1, therefore the determinant of A is also equal to 1. Moreover, the characteristic equation of A shows A always has an eigenvalue of 1 [12]. Thus the eigenvalue of A is equal to 1 with an algebraic multiplicity of 3, implying A is the identity matrix, or the eigenvalues are unique. If we assume A has 3 distinct eigenvalues, because $A \in \mathbb{R}^{3 \times 3}$ and $A = A^T$, we may choose the eigenvectors of A such that they are real. But because $AR = RA$, for every eigenvector, v of A , Rv is also an eigenvector. For an eigenvalue with algebraic multiplicity 1, the corresponding eigenspace is $\dim(1)$, thus $Rv = \mu v$ for some scalar μ , implying each eigenvector of A is also an eigenvector of R . If R is $SO(3)$, however, it has complex eigenvectors in general, which contradicts this assertion. Therefore A must be the identity matrix, and $\hat{K}\hat{K}^T = KK^T$. Referring to K from (2), we observe,

$$KK^T = \begin{bmatrix} FF + X_0X_0^T & X_0 \\ X_0^T & 1 \end{bmatrix}, \quad (8)$$

which implies $\hat{X}_0 = X_0$ and $\hat{F} = F$, or $\hat{K} = K$. ■

It is clear from above that for $R = I$, the relation $AR = RA$ holds trivially, and \hat{K} cannot be uniquely determined. Thus the tolerances with which F in (2) can be determined (in units of pixels) with respect to the amount of camera rotation that occurs is quantified as,

$$\delta f_x < \frac{2f_x^2}{w^2 r_y}; \quad \delta f_y < \frac{2f_y^2}{h^2 r_x}, \quad (9)$$

where r_x and r_y are the x and y -axis rotation angles (in radians) between adjacent frames, and w and h are the width and height of the image, respectively. For a complete proof on the relation between the strength of supervision on K and the closeness of R to I , see [12].

E. Network Architecture

Our framework is inspired by [30]–[32], but rather than requiring two separate networks for depth and pose estimation, we use a common encoder for both tasks in a novel single-network architecture (Fig. 3). That is, given two temporally adjacent input images at times t and t' , our network first convolves these inputs through a series

of convolutional blocks in a common encoder, and then predicts either disparities through a decoder, or camera pose and intrinsics through fully connected layers. In the decoder network, the encoder’s latent representation of the input images is re-upscaled using transposed convolutions with pooling indices from the encoder to fuse low-level features, as inspired by [19], [30]. We use rectified linear units (ReLU) [46] as activation functions in all layers of this decoder except for the prediction layer, which uses a sigmoid function instead. The decoder predicts left-to-right and right-to-left disparities D at both timesteps, which are then either used to reconstruct the right stereo images for a spatially-constrained geometric loss during training via bilinear sampling, or used to construct the depth during inference. In the fully connected layers, translation $\hat{t}_{t' \rightarrow t}$, rotation $\hat{R}_{t' \rightarrow t}$, and camera intrinsics \hat{K} are predicted independently in three separate and decoupled groups of fully connected layers for better performance [31]. These outputs are then either taken at face value during inference as the predicted egomotion and camera parameters, or used as inputs (along with the estimated depth map) to warp the current frame to the next for our temporal reconstruction loss as described previously.

III. RESULTS

In this section, we evaluate our proposed framework using the KITTI driving dataset [47]. Network architecture was implemented using the TensorFlow framework [48] and models were trained on a single NVIDIA GeForce RTX 2070 Super GPU with 8GB of memory using a batch size of 4. Adam optimizer [49] was used to train the network parameters, with exponential decay rates $\beta_1 = 0.9$ and $\beta_2 = 0.99$ and learning rate α initially set to 0.001 but gradually decreased throughout training. Standard data augmentation techniques were used to increase the size of the dataset during training.

We compare our method against the current state-of-the-art using conventional metrics (i.e., Abs Rel, Sq Rel, RMSE, RMSE log, and Accuracy for depth, and Absolute Trajectory Error (ATE) for egomotion) as per [50]. Table I provides a quantitative overview of how our work fits in current literature, and Fig. 4 accompanies this table with a visual, qualitative comparison. Tables II and III evaluate odometry and camera intrinsics estimation performance, and Fig. 5 shows a comparison of trajectory on the KITTI Odometry dataset. We additionally evaluate our framework’s run-time performance to show the benefits of our reduced network size and overhead while maintaining comparable performance with current methods which can be seen in Fig. 6.

A. Performance of Depth Estimation

To evaluate the performance of our network’s depth inference, we use a standard Eigen split [34] on the KITTI dataset [47] as per convention and compare against several state-of-the-art methods with a depth cap of 80m. We train, validate, and test our network using these splits and compare our depth estimation accuracy against multiple other works across several metrics, as shown in Table I. Ground truth data

for the testing set is calculated via projecting the Velodyne LiDAR data onto the image plane.

Several important observations can be extracted from Table I. First, and most notably, our method is the only one which provides a combined architecture for both pose and depth regression (as indicated in the “Comb.” column). Whereas all previous methods split egomotion and depth estimation tasks into separate pipelines, we reduce network redundancy (and hence, number of parameters needed for training) by sharing the input latent representation for both tasks via a common encoder. A second observation is that only two other works before ours are also capable of estimating camera intrinsics during the inference phase ([12], [25], as indicated in the “Int.” column). From these observations, our proposed method, to the best of our knowledge, is the first to demonstrate a simultaneous regression of pose, depth, and camera parameters in a reduced single-network design that uses both spatial and temporal geometric constraints.

However, this reduction in network complexity may come with trade-offs. First, through a quantitative lens, Table I shows that our method is not the lowest in depth estimation error or highest in accuracy. With a reduced network complexity, a possible explanation lies in our encoder’s shared network parameters that must balance both depth and pose/intrinsics pathways. However, even then, our error and accuracy is still strongly comparable to many state-of-the-art methods, many of which have additional components to compensate for occlusions, motion, etc. Compared to the current best with the lowest error and highest accuracy [36], we are on average $\sim 75.9\%$ as error-free and $\sim 95.6\%$ as accurate. This can also be seen qualitatively in Fig. 4. Our method lacks slightly in sharpness compared to [13] but can provide finer edges than [41] and [30]. Depending on one’s setup, the benefits of faster training through a more compact network could outweigh such trade-offs.

B. Learned Camera Intrinsics

To evaluate our system’s ability to recover camera intrinsics (i.e., f_x, f_y, x_0, y_0) through the supervisory signal provided by the rotational component of (1), we follow a similar procedure as [12] and trained separate models on several different video sequences until convergence of these parameters for multiple independent results. During training, parameters were randomly initialized to begin with and empirically had no convergence issues throughout our experiments. We used ten video sequences of the “2011_09_28” subdataset chosen to have the same ground truth calibration done that day, and Table III shows the resulting mean and standard deviation of those ten tests. All experiments were done on the left stereo color camera (“image_02”) of the vehicle setup. For all four variables, we observe that, on average across all ten tests, our method can learn the parameters accurately and within a reasonable bound.

C. Egomotion

We carried out our pose estimation performance evaluation using four sequences from the KITTI Odometry dataset

TABLE I: Comparison of monocular depth estimation with state-of-the-art approaches. Cropped regions from [19] were used for performance evaluation all methods. In the column labeled “Type”, “D” indicates supervised training with ground truth, “T” indicates temporal training only, “S” indicates spatial training only, and “ST” indicates a spatio-temporal training approach. Column “A.C.” denotes whether additional components (such as pre-/post-processing methods) were used in addition to neural networks, and column “Comb.” denotes whether pose and depth networks were combined. Column “Int.” denotes whether that method can simultaneously regress camera intrinsics {Yes (Y), No (N)} (a dash “-” indicates no applicability). We evaluate using the Eigen split [34] on the KITTI dataset [47] and cap depth to 80m as per standard practice [19]. Results from other methods were taken from their corresponding papers. For error metrics, lower is better; for accuracy, higher is better.

Method	Type	A.C.	Comb.	Int.	Error Metrics				Accuracy Metrics		
					Abs Rel	Sq Rel	RMSE	RMSE _{log}	$\delta < 1.25$	$\delta < 1.25^2$	$\delta < 1.25^3$
Train Set Mean	D	-	-	-	0.361	4.826	8.102	0.377	0.638	0.804	0.894
Zou <i>et al.</i> [17]	T	Y	N	N	0.150	1.124	5.507	0.223	0.806	0.933	0.973
Yin <i>et al.</i> [16]	T	Y	N	N	0.149	1.060	5.567	0.226	0.796	0.935	0.975
Chen <i>et al.</i> [25]	T	N	N	Y	0.135	1.070	5.230	0.210	0.841	0.948	0.980
Gordon <i>et al.</i> [12]	T	Y	N	Y	0.128	0.959	5.230	-	-	-	-
Guizilini <i>et al.</i> [36]	T	Y	N	N	0.111	0.785	4.601	0.189	0.878	0.960	0.982
Garg <i>et al.</i> [20]	S	N	-	N	0.177	1.169	5.285	0.282	0.727	0.896	0.958
Godard <i>et al.</i> [19]	S	N	-	N	0.148	1.344	5.927	0.247	0.803	0.922	0.964
Poggi <i>et al.</i> [18]	S	N	-	N	0.163	1.399	6.253	0.262	0.759	0.911	0.961
Pillai <i>et al.</i> [29]	S	N	N	N	0.116	0.935	5.158	0.210	0.842	0.945	0.977
Luo <i>et al.</i> [41]	ST	Y	N	N	0.127	0.936	5.008	0.209	0.841	0.946	0.979
Godard <i>et al.</i> [13]	ST	Y	N	N	0.127	1.031	5.266	0.221	0.836	0.943	0.974
Li <i>et al.</i> [31]	ST	N	N	N	0.183	1.730	6.570	0.268	-	-	-
Babu <i>et al.</i> [32]	ST	N	N	N	0.139	1.174	5.590	0.239	0.812	0.930	0.968
Zhan <i>et al.</i> [30]	ST	N	N	N	0.144	1.391	5.869	0.241	0.803	0.928	0.969
Ours	ST	N	Y	Y	0.141	1.227	5.629	0.239	0.809	0.927	0.962

TABLE II: Comparison of our system’s odometry estimation against various other state-of-the-art methods [15], [32], [51] using absolute trajectory error for translation (*tate*) and rotational (*rate*) movement. Comparison was done on four sequences of the KITTI dataset.

Seq.	Ours		UnDEMoN [32]		SfMLearner [15]		VISO-M [51]	
	<i>tate</i>	<i>rate</i>	<i>tate</i>	<i>rate</i>	<i>tate</i>	<i>rate</i>	<i>tate</i>	<i>rate</i>
00	0.0712	0.0014	0.0644	0.0013	0.7366	0.0040	0.1747	0.0009
04	0.0962	0.0016	0.0974	0.0008	1.5521	0.0027	0.2184	0.0009
05	0.0689	0.0009	0.0696	0.0009	0.7260	0.0036	0.3787	0.0013
07	0.0753	0.0013	0.0742	0.0011	0.5255	0.0036	0.4803	0.0018

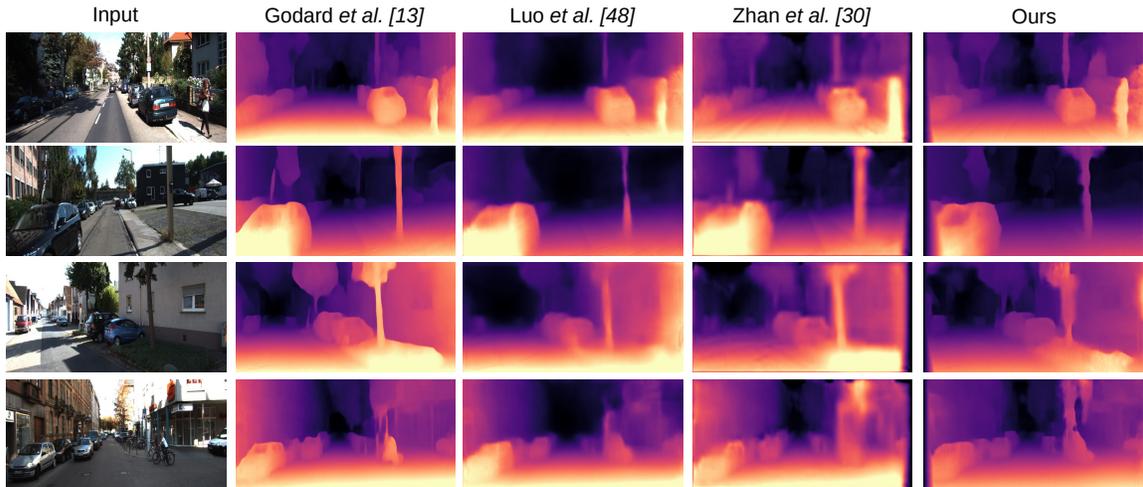


Fig. 4. **Qualitative Comparison of Depths.** Visual comparison of regressed depth maps between our method and various state-of-the-art methods ([13], [30], [41]) on four images from the KITTI Eigen split (images of other methods were retrieved from [13]). Even with a reduced size and complexity, our network can accurately regress depth maps given a single monocular image.

TABLE III: Regressed camera intrinsics during training as compared to the ground truth. Note that ground truth values have been adjusted to match the scaling and cropping done for training. All values are in units of pixels.

Camera Parameter	Learned	Ground Truth
Horizontal Focal Length (f_x)	298.4 ± 2.3	295.8
Vertical Focal Length (f_y)	483.1 ± 3.6	489.2
Horizontal Principal Point (x_0)	254.8 ± 2.4	252.7
Vertical Principal Point (y_0)	127.8 ± 1.7	124.9

[52] and compared against several state-of-the-art methods, including UnDEMoN [32], SfMLearner [15], and VISO-M [51]. For a quantitative comparison, we adopt the absolute trajectory root-mean-square error (ATE) for both translational (t_{ate}) and rotational (r_{ate}) components per standard practice [53]. We note that we used the same model that was trained for depth estimation to output our egomotion estimation, and that these four test sequences were not part of our training set. Sequence 07 is shown in Fig. 5.

From Table II, we observe that for both translational and rotational errors in all four sequences, our method outperformed SfMLearner [15] and VISO-M [51] and is comparable with UnDEMoN’s [32] performance. In contrast to these methods, our system co-predicts egomotion and camera intrinsics (alongside disparity) in a single network such that the loss functions for these free parameters are tied together. This may explain the slight loss in accuracy, especially when compared to [32], but the upside is that our method is a reduction in computational complexity as there are fewer weights in our architecture to optimize over.

D. Run-Time Evaluation

Our single-network architecture design via a common encoder decreases overall network complexity (and therewith the number of network parameters) which can decrease the necessary time to optimize weights. Previously, when using a 7-layer CNN instead of our common encoder for pose and intrinsics regression, network size was ~ 30 million in trainable parameters; however, after replacing those layers with a common encoder, the number of trainable parameters reduced to ~ 27.6 million. This is around an 8% savings.

Fig. 6 shows the effects of this decrease through a run-time comparison. In this figure, the mean (solid) and standard deviation (shaded) loss for each epoch is calculated across ten independent runs for each model. We observed a smaller initial average loss in the combined network, where in this case initialization of weights from the 7-layer CNN that may contribute to a higher loss before being optimal is not necessary. Over time though, we observed that these losses cross paths (roughly at 20 epochs), which is likely caused by the additional 2.4 million parameters for its function approximation. Thus, for those looking to maximize accuracy, a separated network may be better; however, for others who need reasonable results very quickly, a combined network can provide that in just a few epochs.

IV. DISCUSSION

In this work we have presented an unsupervised, single-network monocular depth inference approach for joint pre-

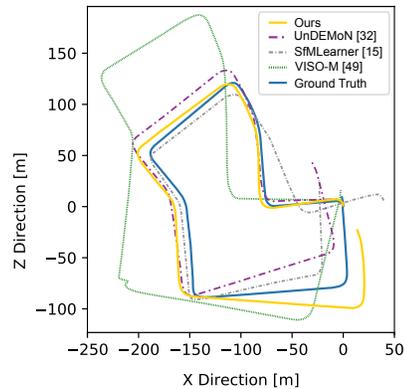


Fig. 5. **Trajectory Comparison.** Visual comparison of estimated egomotion between our method and several others ([15], [32], [51]) on Sequence 07 of the KITTI Odometry dataset. Corresponding t_{ate} and r_{ate} metrics can be found in Table II.

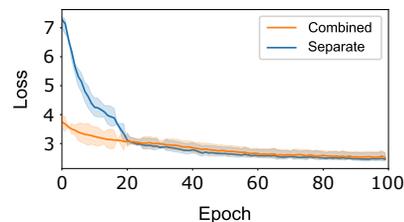


Fig. 6. **Loss Comparison.** Average training loss of the first 100 epochs for both architecture variants. The solid lines represent the mean across ten different runs, and the shaded areas represent one standard deviation ($\sim 70\%$ confidence). The “Separate” architecture used a 7-layer CNN for pose and intrinsics, while “Combined” used a common encoder.

diction of environmental depth, egomotion, and camera intrinsics. Through training our neural network to learn spatial and temporal constraints between stereo and temporally-adjacent pairs, we are able to resolve solutions at metric scale using only monocular video at test time. We distinguish our work from other monocular inference approaches by creating a single, fully differentiable architecture for depth prediction and visual odometry. To further reduce human effort and manual intervention, we also take advantage of intrinsics observability in the system by learning the camera parameters embedded within the temporal reconstruction loss. We verify the success of our system using the KITTI dataset, where our results show we are able to achieve performance comparable to the state-of-the-art in monocular vision while solving for intrinsics and decreasing overhead and overall training complexity.

In future work we plan to quantify network robustness to initialization error during the training of camera parameters. We are also interested in comparing depth and odometry results between predetermined and learned intrinsics, to analyze their effects on prediction outcomes. To improve performance, we expect that the expansion of training to other datasets will provide a more diverse collection of scenes for evaluation, and learned intrinsics will allow for pooling of datasets as another avenue for training. Additionally addressing occlusion and moving objects will ensure added support for higher complexity scenes captured within these datasets.

REFERENCES

- [1] S. Royo and M. Ballesta-Garcia, "An overview of lidar imaging systems for autonomous vehicles," *Applied Sciences*, vol. 9, 2019.
- [2] R. Katzenbeisser, "About the calibration of lidar sensors," in *ISPRS Workshop*, 2003, pp. 1–6.
- [3] N. Muhammad and S. Lacroix, "Calibration of a rotating multi-beam lidar," in *2010 IEEE/RSJ International Conference on Intelligent Robots and Systems*. IEEE, 2010, pp. 5648–5653.
- [4] B. T. Lopez and J. P. How, "Aggressive collision avoidance with limited field-of-view sensing," in *2017 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE, 2017.
- [5] S. M. Seitz, B. Curless, J. Diebel, D. Scharstein, and R. Szeliski, "A comparison and evaluation of multi-view stereo reconstruction algorithms," in *IEEE CVPR*, vol. 1, 2006, pp. 519–528.
- [6] J. L. Schönberger, E. Zheng, J.-M. Frahm, and M. Pollefeys, "Pixel-wise view selection for unstructured multi-view stereo," in *European Conference on Computer Vision*. Springer, 2016, pp. 501–518.
- [7] Y. Furukawa, B. Curless, S. M. Seitz, and R. Szeliski, "Towards internet-scale multi-view stereo," in *2010 IEEE computer society conference on computer vision and pattern recognition*. IEEE, 2010.
- [8] R. Hartley and A. Zisserman, *Multiple view geometry in computer vision*. Cambridge university press, 2003.
- [9] R. A. Newcombe, S. J. Lovegrove, and A. J. Davison, "Dtam: Dense tracking and mapping in real-time," in *2011 ICCV*. IEEE, 2011.
- [10] J. L. Schönberger and J.-M. Frahm, "Structure-from-motion revisited," in *Proceedings of the IEEE CVPR*, 2016, pp. 4104–4113.
- [11] F. Dellaert, S. M. Seitz, C. E. Thorpe, and S. Thrun, "Structure from motion without correspondence," in *Proceedings IEEE Conference on Computer Vision and Pattern Recognition. CVPR 2000 (Cat. No. PR00662)*, vol. 2. IEEE, 2000, pp. 557–564.
- [12] A. Gordon, H. Li, R. Jonschkowski, and A. Angelova, "Depth from videos in the wild: Unsupervised monocular depth learning from unknown cameras," in *IEEE/CVF ICCV*, 2019.
- [13] C. Godard, O. Mac Aodha, M. Firman, and G. Brostow, "Digging Into Self-Supervised Monocular Depth Estimation," Aug. 2019.
- [14] J.-W. Bian, Z. Li, N. Wang, H. Zhan, C. Shen, M.-M. Cheng, and I. Reid, "Unsupervised Scale-consistent Depth and Ego-motion Learning from Monocular Video," *arXiv:1908.10553 [cs]*, Oct. 2019.
- [15] T. Zhou, M. Brown, N. Snavely, and D. G. Lowe, "Unsupervised Learning of Depth and Ego-Motion from Video," in *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. Honolulu, HI: IEEE, Jul. 2017, pp. 6612–6619.
- [16] Z. Yin and J. Shi, "GeoNet: Unsupervised Learning of Dense Depth, Optical Flow and Camera Pose," *arXiv:1803.02276 [cs]*, Mar. 2018.
- [17] Y. Zou, Z. Luo, and J.-B. Huang, "DF-net: Unsupervised joint learning of depth and flow using cross-task consistency," 2018.
- [18] M. Poggi, F. Aleotti, F. Tosi, and S. Mattoccia, "Towards real-time unsupervised monocular depth estimation on CPU," Jul. 2018.
- [19] C. Godard, O. M. Aodha, and G. J. Brostow, "Unsupervised monocular depth estimation with left-right consistency," in *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017.
- [20] R. Garg, V. K. BG, G. Carneiro, and I. Reid, "Unsupervised CNN for Single View Depth Estimation: Geometry to the Rescue," 2016.
- [21] A. Saxena, S. H. Chung, and A. Y. Ng, "Learning depth from single monocular images," in *NeurIPS*, 2006.
- [22] F. Liu, C. Shen, G. Lin, and I. Reid, "Learning depth from single monocular images using deep convolutional neural fields," *IEEE Transactions on PAMI*, vol. 38, no. 10, pp. 2024–2039, 2015.
- [23] A. Saxena, M. Sun, and A. Y. Ng, "Make3d: Learning 3d scene structure from a single still image," *IEEE Transactions on PAMI*, vol. 31, no. 5, pp. 824–840, 2008.
- [24] K. Tateno, F. Tombari, I. Laina, and N. Navab, "CNN-SLAM: Real-time dense monocular SLAM with learned depth prediction," *arXiv:1704.03489 [cs]*, Apr. 2017.
- [25] C. Schmid, C. Sminchisescu, and Y. Chen, "Self-supervised learning with geometric constraints in monocular video - connecting flow, depth, and camera," in *ICCV*, 2019.
- [26] W. N. Greene and N. Roy, "Metrically-scaled monocular slam using learned scale factors," in *2020 IEEE ICRA*, 2020, pp. 43–50.
- [27] H. Fu, M. Gong, C. Wang, K. Batmanghelich, and D. Tao, "Deep ordinal regression network for monocular depth estimation," in *Proceedings of the IEEE CVPR*, 2018, pp. 2002–2011.
- [28] D. Xu, W. Wang *et al.*, "Structured attention guided convolutional neural fields for monocular depth estimation," in *Proceedings of the IEEE CVPR*, 2018.
- [29] S. Pillai, R. Ambrus, and A. Gaidon, "Superdepth: Self-supervised, super-resolved monocular depth estimation," in *2019 International Conference on Robotics and Automation (ICRA)*. IEEE, 2019.
- [30] H. Zhan, R. Garg, C. S. Weerasekera, K. Li, H. Agarwal, and I. Reid, "Unsupervised Learning of Monocular Depth Estimation and Visual Odometry with Deep Feature Reconstruction," Apr. 2018.
- [31] R. Li, S. Wang, Z. Long, and D. Gu, "UnDeepVO: Monocular Visual Odometry through Unsupervised Deep Learning," Feb. 2018.
- [32] V. Madhu Babu, K. Das, A. Majumdar, and S. Kumar, "UnDEMoN: Unsupervised Deep Network for Depth and Ego-Motion Estimation," in *2018 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, Oct. 2018, pp. 1082–1088, iSSN: 2153-0866.
- [33] A. Agha, K. Otsu, B. Morrell, D. D. Fan, R. Thakker, A. Santamaria-Navarro, S.-K. Kim *et al.*, "Nebula: Quest for robotic autonomy in challenging environments; team costar at the darpa subterranean challenge," *arXiv preprint arXiv:2103.11470*, 2021.
- [34] D. Eigen, C. Puhrsch, and R. Fergus, "Depth map prediction from a single image using a multi-scale deep network," in *Advances in Neural Information Processing Systems 27*, Z. Ghahramani, M. Welling, C. Cortes, N. D. Lawrence, and K. Q. Weinberger, Eds. Curran Associates, Inc., 2014, pp. 2366–2374.
- [35] D. Eigen and R. Fergus, "Predicting depth, surface normals and semantic labels with a common multi-scale convolutional architecture,"
- [36] V. Guizilini, R. Ambrus, S. Pillai, A. Raventos, and A. Gaidon, "3d packing for self-supervised monocular depth estimation," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 2485–2494.
- [37] Z. Yang, P. Wang, Y. Wang, W. Xu, and R. Nevatia, "Lego: Learning edge with geometry all at once by watching videos," in *Proceedings of the IEEE CVPR*, 2018, pp. 225–234.
- [38] J. J. Koenderink and A. J. Van Doorn, "Affine structure from motion," *JOSA A*, vol. 8, no. 2, pp. 377–385, 1991.
- [39] R. Clark, S. Wang, H. Wen, A. Markham, and N. Trigoni, "Vinet: Visual-inertial odometry as a sequence-to-sequence learning problem,"
- [40] S. Pillai and J. J. Leonard, "Towards visual ego-motion learning in robots," 2017.
- [41] C. Luo, Z. Yang *et al.*, "Every pixel counts++: Joint learning of geometry and motion with 3d holistic understanding," *IEEE Transactions on PAMI*, vol. 42, pp. 2624–2641, 2019.
- [42] V. Badrinarayanan, A. Kendall, and R. Cipolla, "Segnet: A deep convolutional encoder-decoder architecture for image segmentation," *IEEE Transactions on PAMI*, vol. 39, no. 12, pp. 2481–2495, 2017.
- [43] S. Vijayanarasimhan, S. Ricco *et al.*, "Sfm-net: Learning of structure and motion from video," 2017.
- [44] P. Heise, S. Klose, B. Jensen, and A. Knoll, "Pm-huber: Patchmatch with huber regularization for stereo matching," in *2013 IEEE International Conference on Computer Vision*, 2013, pp. 2360–2367.
- [45] Z. Wang, A. Bovik, H. Sheikh, and E. Simoncelli, "Image Quality Assessment: From Error Visibility to Structural Similarity," *IEEE Transactions on Image Processing*, vol. 13, no. 4, Apr. 2004.
- [46] V. Nair and G. E. Hinton, "Rectified linear units improve restricted boltzmann machines," in *ICML*, 2010.
- [47] M. Menze and A. Geiger, "Object scene flow for autonomous vehicles," in *CVPR*, 2015.
- [48] M. Abadi, A. Agarwal, P. Barham, E. Brevdo, Z. Chen, and *et al.*, "TensorFlow: Large-scale machine learning on heterogeneous systems," 2015. [Online]. Available: <http://tensorflow.org/>
- [49] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," *arXiv preprint arXiv:1412.6980*, 2014.
- [50] C. Zhao, Q. Sun, C. Zhang, Y. Tang, and F. Qian, "Monocular Depth Estimation Based On Deep Learning: An Overview," *Science China Technological Sciences*, vol. 63, no. 9, pp. 1612–1627, Sep. 2020, arXiv: 2003.06620.
- [51] A. Geiger, J. Ziegler, and C. Stiller, "Stereoscan: Dense 3d reconstruction in real-time," in *2011 IEEE intelligent vehicles symposium*. Ieee, 2011, pp. 963–968.
- [52] A. Geiger, P. Lenz, and R. Urtasun, "Are we ready for autonomous driving? the kitti vision benchmark suite," in *2012 IEEE Conference on CVPR*. IEEE, 2012, pp. 3354–3361.
- [53] J. Sturm, N. Engelhard, F. Endres, W. Burgard, and D. Cremers, "A benchmark for the evaluation of rgb-d slam systems," in *2012 IEEE/RSJ International Conference on Intelligent Robots and Systems*. IEEE, 2012, pp. 573–580.