

Unsupervised Depth Completion from Visual Inertial Odometry

Alex Wong[†], Xiaohan Fei[†], Stephanie Tsuei, and Stefano Soatto

Abstract—We describe a method to infer dense depth from camera motion and sparse depth as estimated using a visual-inertial odometry system. Unlike other scenarios using point clouds from lidar or structured light sensors, we have few hundreds to few thousand points, insufficient to inform the topology of the scene. Our method first constructs a piecewise planar scaffolding of the scene, and then uses it to infer dense depth using the image along with the sparse points. We use a predictive cross-modal criterion, akin to “self-supervision,” measuring photometric consistency across time, forward-backward pose consistency, and geometric compatibility with the sparse point cloud. We also present the first visual-inertial + depth dataset, which we hope will foster additional exploration into combining the complementary strengths of visual and inertial sensors. To compare our method to prior work, we adopt the unsupervised KITTI depth completion benchmark, where we achieve state-of-the-art performance. Code available at: <https://github.com/alexklwong/unsupervised-depth-completion-visual-inertial-odometry>.

Index Terms—Visual Learning, Sensor Fusion

I. INTRODUCTION

A sequence of images is a rich source of information about both the three-dimensional (3D) shape of the environment and the motion of the sensor within. Motion can be inferred at most up to a scale and a global Euclidean reference frame, provided sufficient parallax and a number of visually discriminative Lambertian regions that are fixed in the environment and visible from the camera. The position of such regions in the scene defines the Euclidean reference frame, with respect to which motion is estimated. Scale, as well as two directions of orientation, can be further identified by fusion with inertial measurements (accelerometers and gyroscopes) and, if available, a magnetometer can fix the last (Gauge) degree of freedom.

Because the regions defining the reference frame have to be visually distinctive (“features”), they are typically *sparse*. In theory, three points are sufficient to define a Euclidean Gauge if visible at all times. In practice, because of occlusions, any Structure From Motion (SFM) or simultaneous localization and mapping (SLAM) system maintains an estimate of the

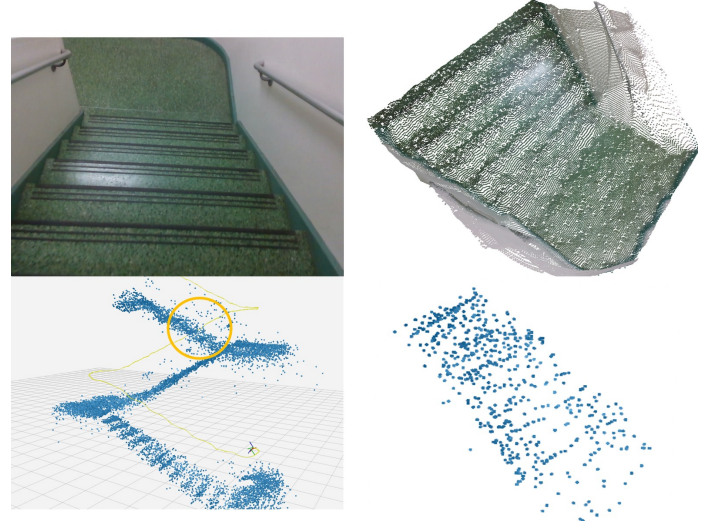


Fig. 1. Depth completion with Visual-Inertial Odometry (VIO) on the proposed VOID dataset (best viewed in color at $5\times$). Bottom left: sparse reconstruction (blue) and camera trajectory (yellow) from VIO. The highlighted region is densified and zoomed in on the top right. Top left shows an image of the same region which is taken as input, and fused with the sparse depth image by our method. On the bottom right is the same view showing only the sparse points, insufficient to determine scene geometry and topology.

location of a sparse set of features, or “sparse point cloud,” typically in the hundreds to thousands. These are sufficient to support a point-estimate of motion, but a rather poor representation of shape as they do not reveal the topology of the scene: The empty space between points could be empty, or occupied by a solid with a smooth surface radiance (appearance). Attempts to *densify* the sparse point cloud, by interpolation or regularization with generic priors such as smoothness, piecewise planarity and the like, typically fail since SFM yields far too sparse a reconstruction to inform topology. This is where the image comes back in.

Inferring shape is ill-posed, even if the point cloud was generated with a lidar or structured light sensor. Filling the gaps relies on assumptions about the environment. Rather than designing ad-hoc priors, we wish to use the image to inform and restrict the set of possible scenes that are compatible with the given sparse points.

Summary of contributions: We use a predictive cross-modal criterion to score dense depth from images and sparse depth. This kind of approach is sometimes referred to as “self-supervised.” Specifically, our method (i) exploits a set of constraints from temporal consistency (a.k.a. photometric consistency across temporally adjacent frames) to pose (forward-backward) consistency in a combination that has not been previously explored. To enable our pose consistency term, we

Manuscript received: September 10, 2019; Revised December 4, 2019; Accepted January 8, 2020.

This paper was recommended for publication by Editor Cesar Cadena upon evaluation of the Associate Editor and Reviewers’ comments. This work was supported by ONR N00014-19-1-2066, and ONR N00014-19-1-2229.

The authors are with the Samueli School of Engineering, Computer Science Department, University of California, Los Angeles, USA. Email: {alexw, feixh, stephanietsuei, soatto}@cs.ucla.edu

[†] denotes authors with equal contributions.

Digital Object Identifier (DOI): see top of this page.

introduce (ii) a set of logarithmic and exponential mapping layers for our network to represent motion using exponential coordinates, which we found to improve reconstruction compared to other parametrizations.

The challenge in using sparse depth as a supervisory (feedback) signal is precisely that it is sparse. Information at the points does not propagate to fill the domain where depth is defined. Some computational mechanism to “diffuse the information” from the sparse points to their neighbors is needed. Our approach proposes (iii) a simple method akin to using a piecewise planar “scaffolding” of the scene, sufficient to transfer the supervisory signal from sparse points to their neighbors. This yields a two-stage approach, where the sparse points are first processed to design the scaffolding (“meshing and interpolation”) and then “refined” using the images as well as priors from the constraints just described.

One additional contribution of our approach is (iv) to introduce the first visual-inertial + depth dataset. The role of inertials is to enable reconstruction in *metric* scale, which is critical for robotic applications. Although scale can be obtained via other sensors, e.g., stereo, lidar, and RGB-D, we note they are not as widely available as monocular cameras with inertials (almost every modern phone has it) and consume more power. Since inertial sensors are now ubiquitous and typically co-located with cameras in many mobile devices from phones to cars, we hope this dataset will foster additional exploration into combining the complementary strengths of visual and inertial sensors.

To evaluate our method, since no other visual-inertial + depth benchmark is available, and to facilitate comparison with similar methods, we adopt the KITTI benchmark, where a Velodyne (lidar) sensor provides sparse points with scale, unlike monocular SFM, but like visual-inertial odometry (VIO). Although the biases in lidar are different from VIO, this can be considered a baseline. Note that we only use the monocular stream of KITTI (not stereo) for fair comparison.

The result is a (v) two-stage approach of scaffolding and refining with a network that contains much fewer parameters than competing methods, yet achieves state-of-the-art performance in the “unsupervised” KITTI benchmark (a misnomer). The supervision in the KITTI benchmark is really fusion from separate sensory channels, combined with ad-hoc interpolation and extrapolation. It is unclear whether the benefit from having such data is outweighed by the biases it induces on the estimate, and in any case such supervision does not scale; hence, we forgo (pseudo) ground truth annotations altogether.

II. RELATED WORK

Supervised Depth Completion minimizes the discrepancy between ground truth depth and depth predicted from an RGB image and sparse depth measurements. Methods focus on network topology [1], [2], [3], optimization [4], [5], [6], and modeling [7], [8]. To handle sparse depth, [1] employed early fusion, where the image and sparse depth are convolved separately and the results concatenated as the input to a ResNet encoder. [9] proposed late fusion via a U-net containing two NASNet encoders for image and sparse depth and jointly learned depth and semantic segmentation, whereas [3] used

ResNet encoders for late fusion. [7] proposed a normalized convolutional layer to propagate sparse depth and used a binary validity map as a confidence measure. [8] proposed an upsampling layer and joint concatenation and convolution to deal with sparse inputs. All these methods require per-pixel ground-truth annotation. What is called “ground truth” in the benchmarks is actually the result of data processing and aggregation of many consecutive frames. We skip such supervision and just infer dense depth by learning the cross-modal fusion from the virtually infinite volume of un-annotated data.

Unsupervised Depth Completion methods, such as [1], [10], [3] predict depth by minimizing the discrepancy between prediction and sparse depth input as well as the photometric error between the input image and its reconstruction from other viewpoints available only during training. [1] used Perspective-n-Point (PnP) [11] and Random Sample Consensus (RANSAC) [12] to align monocular image sequences for their photometric term with a second-order smoothness prior. Yet, [1] does not generalize well to indoor scenes that contains many textureless regions (e.g. walls), where PnP with RANSAC may fail. [10] used a local smoothness term, but instead minimized the photometric error between rectified stereo-pairs where pose is known. [3] also leveraged stereo pairs and a more sophisticated photometric loss [13]. [3] replaced the generic smoothness term with a learned prior to regularize their prediction. To accomplish this, [3] requires a conditional prior network (CPN) that is trained on an *additional* dataset (representative of the depth completion dataset) *in a fully-supervised manner* using ground-truth depth. The CPN does not generalize well outside its training domain (e.g. one cannot use a CPN trained on outdoors scenes to regularize depth predictions for indoors). Hence, [3] is essentially *not* unsupervised and has limited applicability. In contrast, our method is trained on monocular sequences, is *fully unsupervised* and does not use any auxiliary ground-truth supervision. Unlike previous works, our method does not require large networks ([9], [1], [10], [3]) nor any complex network operations ([7], [8]). Moreover, our method outperforms [1], [3] on the unsupervised KITTI depth completion benchmark [2] while using fewer parameters.

Rotation Parameterization. To construct the photometric consistency loss during training, an auxiliary pose network is needed if no camera poses are available. While the translational part of the relative pose can be modeled as $T \in \mathbb{R}^3$, the rotational part belongs to the special orthogonal group $R \in SO(3) \doteq \{R \in \mathbb{R}^{3 \times 3} | R^T R = I, \det(R) = +1\}$ [14], which is represented by a 3×3 matrix. [15] uses quaternions, which require an *additional* norm constraint; this is a soft constraint imposed in the loss function, and thus is not guaranteed. [16], [17], [18] use Euler angles which may result in a non-orthogonal rotation matrix due to rounding error from the multiplication of many sine and cosine terms. We use the exponential map on $SO(3)$ to map the output of the pose network to a rotation matrix. Though theoretically similar, we empirically found that the exponential map is more beneficial than the Euler angles in Sec. VII.

Our contributions are a simple, yet effective two-stage approach resulting in a large reduction in network parameters

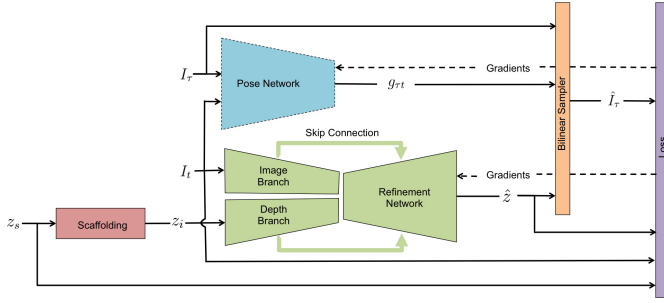


Fig. 2. *System diagram.* (best viewed in color at 5 \times). We first build a scaffolding z_i from sparse depth z_s estimated by VIO. Then together with the image I_t , z_i is fed to the refinement network as input to produce output \hat{z} . Note: the pose network (blue) is only needed in one operation mode and is only used in training. In the other operation mode, VIO poses are used instead. The scaffolding module (red) is parameter-free – leading to our light-weight two-stage approach.

while achieving state-of-the-art performance on the unsupervised KITTI depth completion benchmark; using exponential parameterization of rotation for our pose network; a pose consistency term that enforces forward and backward motion to be the inverse of each other; and finally a new depth completion benchmark for visual-inertial odometry systems with indoor and outdoor scenes and challenging motion.

III. METHOD FORMULATION

We reconstruct a 3D scene given an RGB image $I_t : \mathbb{R}^2 \supset \Omega \mapsto \mathbb{R}_+^3$ and the associated set of sparse depth measurements $z_s : \Omega \supset \Omega_s \mapsto \mathbb{R}_+$.

We begin by assuming that world surfaces are graphs of smooth functions (charts) locally supported on a piecewise planar domain (scaffolding). We construct the scaffolding from the sparse point cloud (“Scaffolding” in Fig. 3) to obtain $z_i : \Omega \mapsto \mathbb{R}_+$, then learn a completion model refining z_i by leveraging the monocular sequences (I_{t-1}, I_t, I_{t+1}) , of frames before and after the given time t , and the sparse depth z_s . We compose a surrogate loss \mathcal{L} (Eqn. 2) for driving the training process, using an encoder-decoder architecture $f_\theta(\cdot)$ parameterized by weights θ , where the input is an image with its scaffolding (I_t, z_i) , and the output is the dense depth $\hat{z} = f_\theta(I_t, z_i)$. Fig. 2 shows an overview of our approach.

A. A Two-Stage Approach

Depth completion is a challenging problem due to the sparsity level of the depth input, z_s . As the density of sparse depth measurements covers $\approx 5\%$ of the image plane for the outdoor self-driving scenario (Sec. V-A) and less than $\approx 1\%$ for the indoor setting (Sec. VII-C), generally only a single measurement will be present within a local neighborhood and in most instances none. This renders *conventional convolutions ineffective* as each sparse depth measurement can be seen as a Dirac delta and convolving a kernel over the entire sparse depth input will give mostly zero activations. Hence, [7], [8], and [2] proposed specialized operations to propagate the information from the sparse depth input through the network. We, instead, propose a two-stage approach that circumvents this problem by first approximating a coarse scene geometry with scaffolding and training a network to refine the approximation.

B. Scaffolding

Given sparse depth measurements z_s , our goal is to create a coarse approximation of the scene; yet, the topology of the scene is not informed by z_s . Hence, we must rely on a prior or an assumption – that surfaces are locally smooth and piecewise planar. We begin by applying the lifting transform [19] to z_s , mapping z_s from 2-d to 3-d space. We then compute its convex hull [20], of which the lower envelope is taken as the Delaunay triangulation of the points in z_s – resulting in a triangular mesh in Barycentric coordinates.

To form the tessellation of the triangular mesh, we approximate each surface using linear interpolation within the Barycentric coordinates and the resulting scaffolding is projected back onto the image plane to produce z_i . For a given triangle, simple interpolation is sufficient for recovering the plane as a linear combination of the co-planar points. For sets of points not co-planar, interpolation will give an approximation, with which we refine using a network. We note that our approximation cannot be achieved by simply filtering (e.g. Gaussian) z_s to propagate depth values as the filter would produce mostly zeros and even destroy the sparse depth signal.

C. Refinement

Given an RGB image and its corresponding piece-wise planar scaffolding (I_t, z_i) , we train a network to recover the 3-d scene by refining z_i based on information from I_t . Our network learns to refine *without* ground-truth supervision by minimizing Eqn. 2 (see Sec. IV).

Network Architecture. We propose two encoder-decoder architectures with skip connections following the late fusion paradigm [9], [3]. Each encoder has an image branch and a depth branch, where each contains 75% and 25% of the total encoder parameters, respectively. The latent representation of the branches are concatenated and fed to the decoder. We propose a VGG11 encoder ($\approx 5.7\text{M}$ parameters) containing 8 convolution layers for each branch as our best performing model, and a VGG8 encoder ($\approx 2.4\text{M}$ parameters) containing only 5 convolution layers for each branch as our light-weight model. This is in contrast to other unsupervised methods [1] (early fusion) and [3] (late fusion) – both of which use ResNet34 encoders with $\approx 23.8\text{M}$ and $\approx 14.8\text{M}$ parameters, respectively. [1], [3] and our approach share the same decoder architecture containing $\approx 4\text{M}$ parameters. We show in Sec. VII that despite having 76.1% and 61.5% fewer encoder parameters than [1] and [3], our VGG11 model outperforms both [1] and [3]. Moreover, performance does not degrade by much from VGG11 to VGG8 and VGG8 still surpasses [1] and [3] while having a 89.9% and 83.9% reduction in the encoder parameters. Unlike [1], [3], which requires high energy consumption hardware, our approach is computationally cheap, and can be deployed to low-powered agents using an Nvidia Jetson.

Logarithmic and Exponential Map Layers. To construct our objective (Eqn. 2), we leverage a pose network [15] to regress the relative camera poses $g = (R, T) \in SE(3) \doteq \{(R, T) | R \in SO(3), T \in \mathbb{R}^3\}$. We present a novel logarithmic map layer: $\log : SO(3) \mapsto so(3)$, where $so(3)$ is the tangent

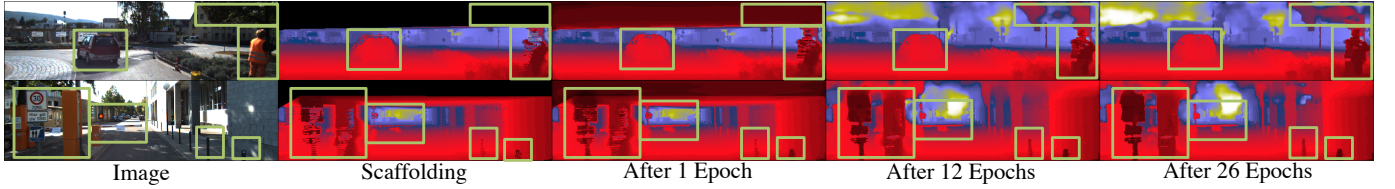


Fig. 3. *Learning to refine* (best viewed at $5\times$ with color). Our network learns to refine the input scaffolding. Green rectangles highlight the regions for comparison throughout the course of training. The network first learns to copy the input and later learns to fuse information from RGB image to refine the approximated depth from scaffolding (see row 1 pedestrian and row 2 street signs).

space of $SO(3)$, and an exponential map layer: $\exp : so(3) \mapsto SO(3)$ – for mapping R between $SO(3)$ and $so(3)$. We use the logarithmic map to construct the pose consistency loss (Eqn. 6), and the exponential to map the output of the pose network $\omega \doteq [\omega_1, \omega_2, \omega_3]^T \in \mathbb{R}^3$ as coordinates in $so(3)$ to a rotation matrix:

$$R(\omega) = \exp(\hat{\omega}) \doteq I + \hat{\omega} \sin \|\omega\|_2 + \hat{\omega}^2 (1 - \cos \|\omega\|_2) \quad (1)$$

where the hat operator $\hat{\cdot}$ maps $\omega \in \mathbb{R}^3$ to a skew-symmetric matrix [14]. We train our pose network using a surrogate loss (Eqn. 3) without explicit supervision. Ablation studies on the use of exponential coordinates and pose consistency for depth completion can be found in Table III and IV.

Our approach contains two stages: (i) we generate a coarse piecewise planar approximation of the scene from the sparse depth inputs z_s via scaffolding and (ii) we feed the resulting depth map along with the associated RGB image to our network for refinement (Fig. 3). This approach *alleviates* the network from the need of learning from sparse inputs, for which [1] and [3] compensated with parameters. We show the effectiveness of this approach by achieving the state-of-the-art on the unsupervised KITTI depth completion benchmark with half as many parameters as the prior-art.

IV. LOSS FUNCTION

Our loss function is a linear combination of four terms that constrain (i) the photometric consistency between the observed image and its reconstructions from the monocular sequence, (ii) the predicted depth to be similar to that of the associated available sparse depth, (iii) the composition of the predicted forward and backward relative poses to be the identity, and (iv) the prediction to adhere to local smoothness.

$$\mathcal{L} = w_{ph} L_{ph} + w_{sz} L_{sz} + w_{pc} L_{pc} + w_{sm} L_{sm} \quad (2)$$

where L_{ph} denotes photometric consistency, L_{sz} sparse depth consistency, L_{pc} pose consistency, and L_{sm} local smoothness. Each loss term L is described in the next subsections and the associated weight w in Sec. VI.

A. Photometric Consistency

We enforce temporal consistency by minimizing the discrepancy between each observed image I_t and its reconstruction \hat{I}_τ from temporally adjacent images I_τ , where $\tau \in T \doteq \{t-1, t+1\}$:

$$\hat{I}_\tau(x) = I_\tau(\pi g_{\tau t} K^{-1} \bar{x} \hat{z}(x)) \quad (3)$$

where $\bar{x} = [x^T 1]^T$ are the homogeneous coordinates of $x \in \Omega$, $g_{\tau t} \in SE(3)$ is the relative pose of the camera from time

t to τ , K denotes the camera intrinsics, and π refers to the perspective projection.

Our photometric consistency term is a combination of the average per pixel reprojection residual with an L_1 penalty and SSIM [13], a perceptual metric that is invariant to local illumination changes:

$$L_{ph} = \frac{1}{|\Omega|} \sum_{\tau \in T} \sum_{x \in \Omega} w_{co} |I_t(x) - \hat{I}_\tau(x)| + w_{st} (1 - \text{SSIM}(I_t(x), \hat{I}_\tau(x))) \quad (4)$$

We use 3×3 image patches centered at location x for SSIM. w_{co} and w_{st} can be found in Sec. VI.

B. Sparse Depth Consistency

Our sparse depth consistency term provides our predictions with *metric* scale by encouraging the predictions \hat{z} to be similar to that of the *metric* sparse depth z_s available from lidar in KITTI dataset (Sec. V-A) and sparse reconstruction in our visual-inertial dataset (Sec. V-B). Our sparse depth consistency loss is the L_1 -norm of the difference between the predicted depth \hat{z} and the sparse depth z_s averaged over Ω_s (the support of the sparse depth):

$$L_{sz} = \frac{1}{|\Omega_s|} \sum_{x \in \Omega_s} |\hat{z}(x) - z_s(x)| \quad (5)$$

C. Pose Consistency

A pose network takes an ordered pair of images (I_t, I_τ) and outputs the relative pose $g_{\tau t} \in SE(3)$ (forward pose). When a temporally swapped pair (I_τ, I_t) is fed to the network, the network is expected to output $g_{t\tau}$ (backward pose) – the inverse of $g_{\tau t}$, i.e., $g_{\tau t} \cdot g_{t\tau} = e \in SE(3)$. The forward-backward pose consistency thus penalizes the deviation of the composed pose from the identity:

$$L_{pc} = \|\log(g_{\tau t} \cdot g_{t\tau})\|_2^2 \quad (6)$$

where $\log : SE(3) \mapsto se(3)$ is the logarithmic map.

D. Local Smoothness

We impose a smoothness loss on the predicted depth \hat{z} by applying an L_1 penalty to the gradients in both the x and y directions of the predicted depth \hat{z} :

$$L_{sm} = \frac{1}{|\Omega|} \sum_{x \in \Omega} \lambda_X(x) |\partial_X \hat{z}(x)| + \lambda_Y(x) |\partial_Y \hat{z}(x)| \quad (7)$$

where $\lambda_X = e^{-|\partial_X I_t(x)|}$ and $\lambda_Y = e^{-|\partial_Y I_t(x)|}$ are the edge-awareness weights to allow for discontinuities in regions corresponding to object boundaries.

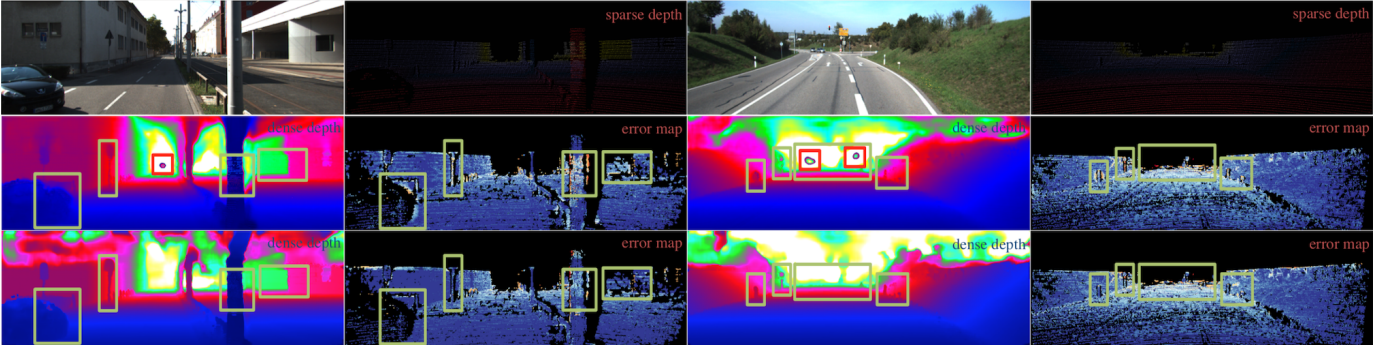


Fig. 4. *Qualitative evaluation on KITTI benchmark.* Top to bottom: input image and sparse depth, results of [1], our results. Results are taken from KITTI online test server. Warmer colors in the error map denote higher error. Green rectangles highlight regions for detail comparison. We perform better in general, particularly on thin structures and far regions. [1] exhibit artifacts resembling scanlines and “circles” for far away regions (highlighted in red).

E. The Role of Inertials

Although inertials are not directly present in the loss, their role in *metric* depth completion is crucial. Without inertials, a SLAM system cannot produce sparse point clouds in metric scale, which are then used as both the input to the scaffolding stage (Sec. III-B) and a supervisory signal (Eqn. 5).

V. DATASETS

A. KITTI Benchmark

We evaluate our approach on the KITTI depth completion benchmark [2]. The dataset provides $\approx 80,000$ raw image frames and associated sparse depth maps. The sparse depth maps are the raw output from the Velodyne lidar sensor, each with a density of $\approx 5\%$. The ground-truth depth map is created by accumulating the neighbouring 11 raw lidar scans, with dense depth corresponding to the bottom 30% of the images. We use the officially selected 1,000 samples for validation and we apply our method to 1,000 testing samples, with which we submit to the official KITTI website for evaluation. The results are reported in Table II.

B. VOID Benchmark

While KITTI provides a standard benchmark for evaluating depth completion in the driving scenario, there exists no standard depth completion benchmark for the indoor scenario. [1], [3] used NYUv2 [21] – an RGB-D dataset – to develop and evaluate their models on indoor scenes. Yet, each performs a different evaluation protocol with different sparse depth samples – varying densities of depth values were randomly sampled from the depth frame, preventing direct comparisons between methods. Though this is reasonable as a proof of concept, it is not realistic in the sense that no sensor measures depth at random locations.

The VOID dataset. We propose a new publicly available dataset for a real world use case of depth completion by bootstrapping sparse reconstruction in *metric* space from a SLAM system. While it is well known that metric scale is not observable in the purely image-based SLAM and SFM setting, it has been resolved by the recent advances in VIO [22], [23], where metric pose and structure estimation can be realized in a gravity-aligned and scaled reference frame using an inertial measurement unit (IMU). To this end, we leverage an off-the-shelf VIO system¹, atop which we construct our dataset and

develop our depth completion model. While there are some visual-inertial datasets (e.g. TUM-VI [24] and PennCOSYVIO [25]), they lack per-frame dense depth measurements for cross-modal validation, and are also relatively small – rendering them unsuitable for training deep learning models. To demonstrate the applicability of our approach, we additionally show qualitative results on the TUM-VI dataset in Fig. 5 using sparse depth density level of 0.015%.

Our dataset is dubbed “Visual Odometry with Inertial and Depth” or “VOID” for short and is comprised of RGB video streams and inertial measurements for *metric* reconstruction along with per-frame dense depth for cross-modal validation. **Data acquisition.** Our data was collected using the latest Intel RealSense D435i camera², which was configured to produce synchronized accelerometer and gyroscope measurements at 400 Hz, along with synchronized VGA-size (640×480) RGB and depth streams at 30 Hz. The depth frames are acquired using active stereo and is aligned to the RGB frame using the sensor factory calibration (see Fig. 8). All the measurements are time-stamped.

The SLAM system we use is based on [22] – an EKF-based VIO model. While the VIO recursively estimates a joint posterior of the state of the sensor platform (e.g. pose, velocity, sensor biases, and camera-to-IMU alignment) and a small set of reliable feature points, the 3D structure it estimates is extremely sparse – typically 20 \sim 30 feature points (in-state features). To facilitate 3D reconstruction, we track a moderate amount of out-of-state features in addition to the in-state ones, and estimate the depth of the feature points using auxiliary depth sub-filters [14].

The benchmark. We evaluate our method on the VOID depth completion benchmark, which contains 56 sequences in total, both indoor and outdoor with challenging motion. Typical scenes include classrooms, offices, stairwells, laboratories, and gardens. Of the 56 sequences, 48 sequences ($\sim 40K$ frames) are designated for training and 8 sequences for testing, from which we sampled 800 frames to construct the testing set. Our benchmark provides sparse depth maps at three density levels. We configured our SLAM system to track and estimate depth of 1500, 500 and 150 feature points, corresponding to 0.5%, 0.15% and 0.05% density of VGA size, which are then used in the depth completion task.

¹<https://github.com/ucla-vision/xivo>

²<https://realsense.intel.com/depth-camera/>

TABLE I
ERROR METRICS.

Metric	units	Definition
MAE	mm	$\frac{1}{ \Omega } \sum_{x \in \Omega} \hat{z}(x) - z_{gt}(x) $
RMSE	mm	$(\frac{1}{ \Omega } \sum_{x \in \Omega} \hat{z}(x) - z_{gt}(x) ^2)^{1/2}$
iMAE	1/km	$\frac{1}{ \Omega } \sum_{x \in \Omega} 1/\hat{z}(x) - 1/z_{gt}(x) $
iRMSE	1/km	$(\frac{1}{ \Omega } \sum_{x \in \Omega} 1/\hat{z}(x) - 1/z_{gt}(x) ^2)^{1/2}$

Error metrics for evaluating KITTI and VOID depth completion benchmarks, where z_{gt} is the ground truth.

VI. IMPLEMENTATION DETAILS

Our approach was implemented using TensorFlow [26]. With a Nvidia GTX 1080Ti, training takes ≈ 42 hours for our VGG11 model and ≈ 34 hours for our VGG8 model on KITTI depth completion benchmark (Sec. V-A) for 30 epochs; whereas training takes ≈ 10 hours and ≈ 7 hours on the VOID benchmark (Sec. V-B) for 10 epochs. Inference takes ≈ 22 ms per image. We used Adam [27] with $\beta_1 = 0.9$ and $\beta_2 = 0.999$ to optimize our network end-to-end with a base learning rates of 1.2×10^{-4} for KITTI and 1×10^{-4} for VOID. We decrease the learning rate by half after 18 epochs for KITTI and 6 epochs for VOID, and again after 24 epochs and 8 epochs, respectively. We train our network with a batch size of 8 using a 768×320 resolution for KITTI and 640×480 for VOID. We are able to achieve our results on the KITTI benchmark using the following set of weights for each term in our loss function: $w_{ph} = 1.00$, $w_{co} = 0.20$, $w_{st} = 0.40$, $w_{sz} = 0.20$, $w_{pc} = 0.10$ and $w_{sm} = 0.01$. For the VOID benchmark, we increased w_{sz} to 1.00 and w_{sm} to 0.10. We do not use any data augmentation.

VII. EXPERIMENTS AND RESULTS

A. KITTI Depth Completion Benchmark

We show quantitative and qualitative comparisons on the unsupervised KITTI depth completion benchmark in Table II and Fig. 4, respectively. The results of the methods listed are taken directly from their papers. We note that [3] only reported their result in their paper and do not have an entry in KITTI depth completion benchmark for their unsupervised model. Hence, we compare qualitatively with the prior-art [1]. Our VGG11 model outperforms the state-of-the-art [3] on every metric by as much as 12.8% while using 48.4% fewer parameters. Our light-weight VGG8 model also outperforms [3] on MAE, RMSE, and iMAE while [3] beat our VGG8 by 2.2% on iRMSE. We note that [3] trains a separate network, using ground truth, to supervise their depth completion model. Moreover, [3] exploits rectified stereo-imagery where the pose of the cameras is known; whereas, we learn our pose by jointly training the pose network with our depth predictor. In comparison to [1] (who also uses monocular videos), both our VGG11 and VGG8 model outperforms them on every metric while using much fewer parameters. We also note that the qualitative results of [1] contains artifacts such as apparent scanlines of the Velodyne and “circles” in far regions.

As an introspective exercise, we plot the mean error of our model at varying distances on the KITTI validation set (Fig. 6) and overlay it with the ground truth depth distribution to show

TABLE II
KITTI DEPTH COMPLETION BENCHMARK.

Method	# Parameters	MAE	RMSE	iMAE	iRMSE
Schneider [28]	not reported	605.47	2312.57	2.05	7.38
Ma [1]	$\approx 27.8M$	350.32	1299.85	1.57	4.07
Yang [3]	$\approx 18.8M$	343.46	1263.19	1.32	3.58
Ours VGG11	$\approx 9.7M$	299.41	1169.97	1.20	3.56
Ours VGG8	$\approx 6.4M$	304.57	1164.58	1.28	3.66

We compare our model to unsupervised methods on the KITTI depth completion benchmark [2]. Number of parameters used by each are listed for comparison. [28] stated that they use a fully convolution network, but does not specify the full architecture. Our VGG11 model outperforms state-of-the-art [3] across all metrics. Despite reducing $\approx 3.3M$ parameters, our VGG8 model does not degrade by much and outperforms VGG11 marginally on the RMSE metric. Moreover, our VGG8 model also outperforms [1] and [3].

TABLE III
KITTI DEPTH COMPLETION ABLATION STUDY.

Model	Encoder	Rot.	MAE	RMSE	iMAE	iRMSE
Scaffolding	-	-	443.57	1990.68	1.72	6.43
$L_{ph} + L_{sz} + L_{sm}$ (vanilla)	VGG11	Eul.	347.14	1330.88	1.46	4.22
$L_{ph} + L_{sz} + L_{sm}$	VGG11	Eul.	327.84	1262.46	1.31	3.87
$L_{ph} + L_{sz} + L_{sm}$	VGG11	Exp.	312.10	1255.21	1.28	3.86
$L_{ph} + L_{sz} + L_{pc} + L_{sm}$	VGG11	Exp.	305.06	1239.06	1.21	3.71
$L_{ph} + L_{sz} + L_{pc} + L_{sm}$	VGG8	Exp.	308.81	1230.85	1.29	3.84

We compare variants of our model on the KITTI depth completion validation set. Each model is denoted by its loss function. Regions with missing depth in Scaffolding Only is assigned average depth. It is clear that scaffolding alone (row 1) and our baseline model trained *without* scaffolding (row 2) do poorly compared to our models that combine both (rows 3-6). Our full model using VGG11 produces the best overall results and achieves state-of-the-art on the test set Table II. Our approach is robust, our light-weight VGG8 model achieves similar performance to our VGG11 model.

that our model performs very well in distances that matter in real-life scenarios. Our performance begins to degrade at distances larger than 80 meters; this is due to the lack of sparse measurements and insufficient parallax – problems that plague methods relying on multi-view supervision.

B. KITTI Depth Completion Ablation Study

We analyze the effect brought by each of our contributions through a quantitative evaluation on the KITTI depth completion validation set (Table III). Our two baseline models, scaffolding and vanilla model trained without scaffolding, perform poorly in comparison to the models that are trained with scaffolding – showcasing the effectiveness of our refinement approach. Although the loss functions are identical, exponential parameterization consistently improves over Euler angles across all metrics. We believe this is due to the regularity of the derivatives of the exponential map [29] compared to other parameterizations – resulting in faster convergence and wider minima during training. While [16], [30], [17] train their pose network using the photometric error with no additional constraint, we show that it is beneficial to impose our pose consistency term (Sec. 6). By constraining the forward and backward poses to be inverse of each other, we obtain a more accurate pose resulting in better depth prediction. Our experiments verify this claim as we see an improvement in across all metrics in Table III. We note that the improvement does not seem significant on KITTI as the motion is mostly planar; however, when predicting non-trivial 6 DoF motion (Sec. VII-D), we see a significant boost when employing this term. Our model trained with the full loss function produces

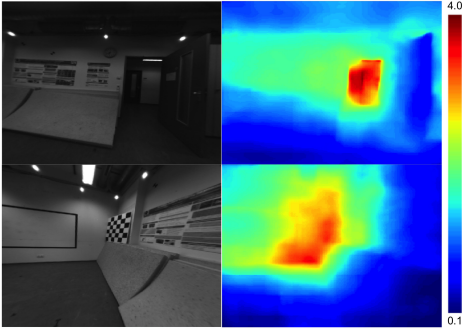


Fig. 5. *Qualitative results on TUM-VI* (best viewed in color at $2\times$). We apply our method to TUM-VI and obtained our results using sparse depth input at a density level of 0.015%. Unlike KITTI and VOID, TUM-VI images are monochrome, and bear a highly distorted fisheye camera model, which was compensated in training. Color bar shows the depth range.

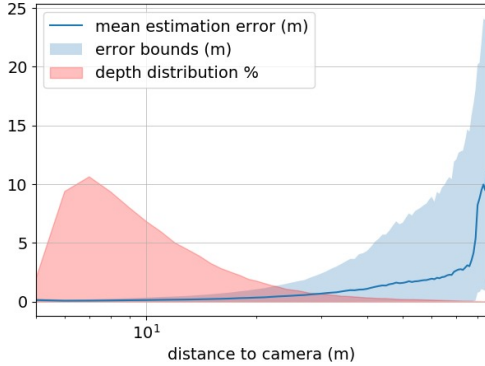


Fig. 6. *Error characteristics of our model on KITTI*. The abscissa shows the distance of sparse data points measured by Velodyne, of which the percentage of all the data points is shown in red; the blue curve shows the mean absolute error of the estimated depth at the given distance, of which the 5-th and 95-th percentile enclose the light blue region.

the best results (bolded in Table II) and is the state-of-the-art for unsupervised KITTI depth completion benchmark. We further propose a VGG8 model that only contains $\approx 6.4\text{M}$ parameters. Despite having 34% fewer parameters than VGG11, the performance of VGG8 does not degrade by much (see Table II, III, V).

C. VOID Depth Completion Benchmark

We evaluate our method on the VOID depth completion benchmark for all three density levels (Table V) using error metrics in Table I. As the photometric loss (Eqn. 4) is largely dependent on obtaining the correct pose, we additionally propose a hybrid model, where the relative camera poses from our visual-inertial SLAM system are used to construct the photometric loss to show an upper bound on performance. In contrast to the KITTI, which provides $\approx 5\%$ sparse depth density concentrated on the bottom 30% of the image, the VOID benchmark only provides $\approx 0.5\%$, $\approx 0.15\%$ and $\approx 0.05\%$ densities in sparse depth. Yet, our method is still able to produce reasonable results for indoor scenes with a MAE of $\approx 8.5\text{ cm}$ on 0.5% density and $\approx 17.9\text{ cm}$ when given only 0.05%. Since most scenes contain textureless regions, sparse depth supervision becomes important as photometric reconstruction is unreliable. Hence, performance degrades as density decreases. Yet, we degrade gracefully: as density decreases by 10X, our error only doubles. We note that the scaffolding may poorly represent the scene. In the worst case,



Fig. 7. *Qualitative evaluation on VOID benchmark*. Top: Input RGB images. Bottom: Densified depth images back-projected to 3D, colored, and viewed from a different vantage point.

where it provides no extra information, our method becomes the common depth completion approach. Also, we observe systematic performance improvement in all the evaluation metrics (Table V) when replacing the pose network with SLAM pose. This can be largely attributed to the necessity for the correct pose to minimize photometric error during training. Our pose network may not be able to consistently predict the correct pose due to the challenging motion of the dataset. Fig. 7 shows two sample RGB images with the densified depth images back-projected to 3D, colored, and viewed from a different vantage point.

D. VOID Depth Completion Ablation Study

To better understand the effect of rotation parameterization and our pose consistency loss (Eqn. 6) on the depth completion task, we compare variants of our model and again replace the pose network with SLAM pose to show an upper-bound on performance. Although exponential outperforms Euler parameterization, we note that both perform much worse than using SLAM pose. However, we observe a performance boost when applying our pose consistency term and our model improves over exponential without pose consistency by as much as 23.4%. Moreover, it approaches the performance of our model trained using SLAM pose. This trend still holds when density decreases (Table V). This suggests that despite the additional constraint, the pose network still has some difficulties predicting the pose due to the challenging motion. This finding, along with results from Table V, highlights the strength of classical SLAM systems in the deep learning era, which also urges us to develop and test pose networks on the VOID dataset which features non-trivial 6 DoF motion – much more challenging than the mostly-planar motion in KITTI.

VIII. DISCUSSION

While deep networks have attracted a lot of attention as a general framework to solve an array of problems, we must note that pose may be difficult to learn on datasets with non-trivial 6 DoF motion – which the SLAM community has studied for decades. We hope that VOID will serve as a platform to develop models that can handle challenging motion and further foster fusion of multi-sensory data. Furthermore, we show that a network can recover the scene geometry from extremely

TABLE IV
VOID DEPTH COMPLETION BENCHMARK AND ABLATION STUDY.

Method	MAE	RMSE	iMAE	iRMSE
Ma [1]	198.76	260.67	88.07	114.96
Yang [3]	151.86	222.36	74.59	112.36
VGG11 PoseNet + Eul.	108.97	212.16	64.54	142.64
VGG11 PoseNet + Exp.	103.31	179.05	63.88	131.06
VGG11 PoseNet + Exp. + L_{pc}	85.05	169.79	48.92	104.02
VGG11 SLAM Pose	73.14	146.40	42.55	93.16
VGG8 PoseNet + Exp. + L_{pc}	94.33	168.92	56.01	111.54

We compare the variants of our pose network. SLAM Pose replaces the output of pose network with SLAM estimated pose to gauge an upper bound in performance. When using our pose consistency term with exponential parameterization, our method approaches the performance of our method when using SLAM pose. Note: we trained [1] from scratch using ground-truth pose and adapted [26] to train on monocular sequences. The conditional prior network used in [3] is trained on ground truth from NYUv2 [21].

TABLE V
DEPTH COMPLETION ON VOID WITH VARYING SPARSE DEPTH DENSITY.

Density	Pose From	MAE	RMSE	iMAE	iRMSE
~ 0.5%	PoseNet	85.05	169.79	48.92	104.02
	SLAM	73.14	146.40	42.55	93.16
~ 0.15%	PoseNet	124.11	217.43	66.95	121.23
	SLAM	118.01	195.32	59.29	101.72
~ 0.05%	PoseNet	179.66	281.09	95.27	151.66
	SLAM	174.04	253.14	87.39	126.30

The VOID dataset contains VGA size images (480×640) of both indoor and outdoor scenes with challenging motion. For “Pose From”, SLAM refers to relative poses estimated by a SLAM system, and PoseNet refers to relative poses predicted by a pose network.

sparse point clouds (e.g. features tracked by SLAM). We also show that improvements can be obtained by leveraging pose from a SLAM system instead of a pose network. These findings motivate a possible mutually beneficial marriage between classical methods and deep learning.

REFERENCES

- [1] F. Ma, G. V. Cavalheiro, and S. Karaman, “Self-supervised sparse-to-dense: Self-supervised depth completion from lidar and monocular camera,” in *International Conference on Robotics and Automation (ICRA)*. IEEE, 2019.
- [2] J. Uhrig, N. Schneider, L. Schneider, U. Franke, T. Brox, and A. Geiger, “Sparsity invariant cnns,” in *2017 International Conference on 3D Vision (3DV)*. IEEE, 2017, pp. 11–20.
- [3] Y. Yang, A. Wong, and S. Soatto, “Dense depth posterior (ddp) from single image and sparse range,” *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019.
- [4] N. Chodosh, C. Wang, and S. Lucey, “Deep Convolutional Compressed Sensing for LiDAR Depth Completion,” in *Asian Conference on Computer Vision (ACCV)*, 2018.
- [5] M. Dimitrievski, P. Veelaert, and W. Philips, “Learning morphological operators for depth completion,” in *Advanced Concepts for Intelligent Vision Systems*, 2018.
- [6] Y. Zhang and T. Funkhouser, “Deep depth completion of a single rgb-d image,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 175–185.
- [7] A. Eldesokey, M. Felsberg, and F. S. Khan, “Propagating confidences through cnns for sparse data regression,” in *Proceedings of British Machine Vision Conference (BMVC)*, 2018.
- [8] Z. Huang, J. Fan, S. Cheng, S. Yi, X. Wang, and H. Li, “Hms-net: Hierarchical multi-scale sparsity-invariant network for sparse depth completion,” *IEEE Transactions on Image Processing*, 2019.
- [9] M. Jaritz, R. de Charette, E. Wirbel, X. Perrotton, and F. Nashashibi, “Sparse and dense data with cnns: Depth completion and semantic segmentation,” in *International Conference on 3D Vision (3DV)*, 2018.
- [10] S. S. Shivakumar, T. Nguyen, I. D. Miller, S. W. Chen, V. Kumar, and C. J. Taylor, “Dfusenet: Deep fusion of rgb and sparse depth information for image guided dense depth completion,” in *2019 IEEE Intelligent Transportation Systems Conference (ITSC)*. IEEE, 2019, pp. 13–20.

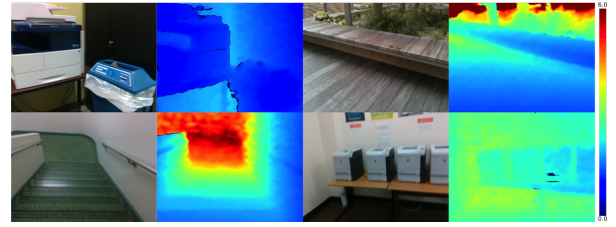


Fig. 8. Sample RGB + D images in the VOID dataset (best viewed in color at $5\times$). Color bar shows the depth range.

- [11] V. Lepetit, F. Moreno-Noguer, and P. Fua, “Epnnp: An accurate o (n) solution to the pnp problem,” *International journal of computer vision*, vol. 81, no. 2, p. 155, 2009.
- [12] M. A. Fischler and R. C. Bolles, “Random sample consensus: a paradigm for model fitting with applications to image analysis and automated cartography,” *Communications of the ACM*, vol. 24, no. 6, 1981.
- [13] Z. Wang, A. C. Bovik, H. R. Sheikh, E. P. Simoncelli, et al., “Image quality assessment: from error visibility to structural similarity,” *IEEE transactions on image processing*, vol. 13, no. 4, pp. 600–612, 2004.
- [14] Y. Ma, S. Soatto, J. Kosecka, and S. Sastry, *An invitation to 3-d vision: from images to geometric models*. Springer Science & Business Media, 2012, vol. 26.
- [15] A. Kendall, M. Grimes, and R. Cipolla, “Posenet: A convolutional network for real-time 6-dof camera relocalization,” in *Proceedings of the IEEE international conference on computer vision*, 2015.
- [16] X. Fei, A. Wong, and S. Soatto, “Geo-supervised visual depth prediction,” *IEEE Robotics and Automation Letters*, vol. 4, no. 2, 2019.
- [17] Z. Yin and J. Shi, “Geonet: Unsupervised learning of dense depth, optical flow and camera pose,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 1983–1992.
- [18] T. Zhou, M. Brown, N. Snavely, and D. G. Lowe, “Unsupervised learning of depth and ego-motion from video,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017.
- [19] K. Q. Brown, “Voronoi diagrams from convex hulls,” *Information processing letters*, vol. 9, no. 5, pp. 223–228, 1979.
- [20] C. B. Barber, D. P. Dobkin, D. P. Dobkin, and H. Huhdanpaa, “The quickhull algorithm for convex hulls,” *ACM Transactions on Mathematical Software (TOMS)*, vol. 22, no. 4, pp. 469–483, 1996.
- [21] P. K. Nathan Silberman, Derek Hoiem and R. Fergus, “Indoor segmentation and support inference from rgb-d images,” in *ECCV*, 2012.
- [22] E. Jones and S. Soatto, “Visual-inertial navigation, mapping and localization: A scalable real-time causal approach,” *International Journal of Robotics Research*, January 2011.
- [23] A. I. Mourikis and S. I. Roumeliotis, “A multi-state constraint kalman filter for vision-aided inertial navigation,” in *Robotics and automation, 2007 IEEE international conference on*. IEEE, 2007, pp. 3565–3572.
- [24] D. Schubert, T. Goll, N. Demmel, V. Usenko, J. Stückler, and D. Cremers, “The tum vi benchmark for evaluating visual-inertial odometry,” in *2018 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE, 2018, pp. 1680–1687.
- [25] B. Pfrommer, N. Sanket, K. Daniilidis, and J. Cleveland, “Pencosyvio: A challenging visual inertial odometry benchmark,” in *2017 IEEE International Conference on Robotics and Automation, ICRA 2017, Singapore, Singapore, May 29 - June 3, 2017*, 2017, pp. 3847–3854.
- [26] M. Abadi, P. Barham, J. Chen, Z. Chen, A. Davis, J. Dean, M. Devin, S. Ghemawat, G. Irving, M. Isard, et al., “Tensorflow: A system for large-scale machine learning,” in *OSDI*, vol. 16, 2016, pp. 265–283.
- [27] D. P. Kingma and J. Ba, “Adam: A method for stochastic optimization,” *arXiv preprint arXiv:1412.6980*, 2014.
- [28] N. Schneider, L. Schneider, P. Pinggera, U. Franke, M. Pollefeys, and C. Stiller, “Semantically guided depth upsampling,” in *German Conference on Pattern Recognition*. Springer, 2016.
- [29] G. Gallego and A. Yezzi, “A compact formula for the derivative of a 3-d rotation in exponential coordinates,” *Journal of Mathematical Imaging and Vision*, vol. 51, no. 3, pp. 378–384, 2015.
- [30] C. Wang, J. Miguel Buenaposada, R. Zhu, and S. Lucey, “Learning depth from monocular videos using direct methods,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 2022–2030.
- [31] 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.
- [32] K. Simonyan and A. Zisserman, “Very deep convolutional networks for large-scale image recognition,” *arXiv preprint arXiv:1409.1556*, 2014.

APPENDIX A VOID DATASET

In the main paper, we introduced the “Visual Odometry with Inertial and Depth” (VOID) dataset with which we propose a new depth completion benchmark. We described the data acquisition process, benchmark setup, and evaluation protocols in Sec. V-B and Sec. VII-C. To give some flavor of the VOID dataset, Fig. 9 shows a set of images (top inset) sampled from video sequences in VOID, and output of our visual-inertial odometry (VIO) system (bottom), where the blue pointcloud is the sparse reconstruction of the underlying scene and the yellow trace is the estimated camera trajectory.

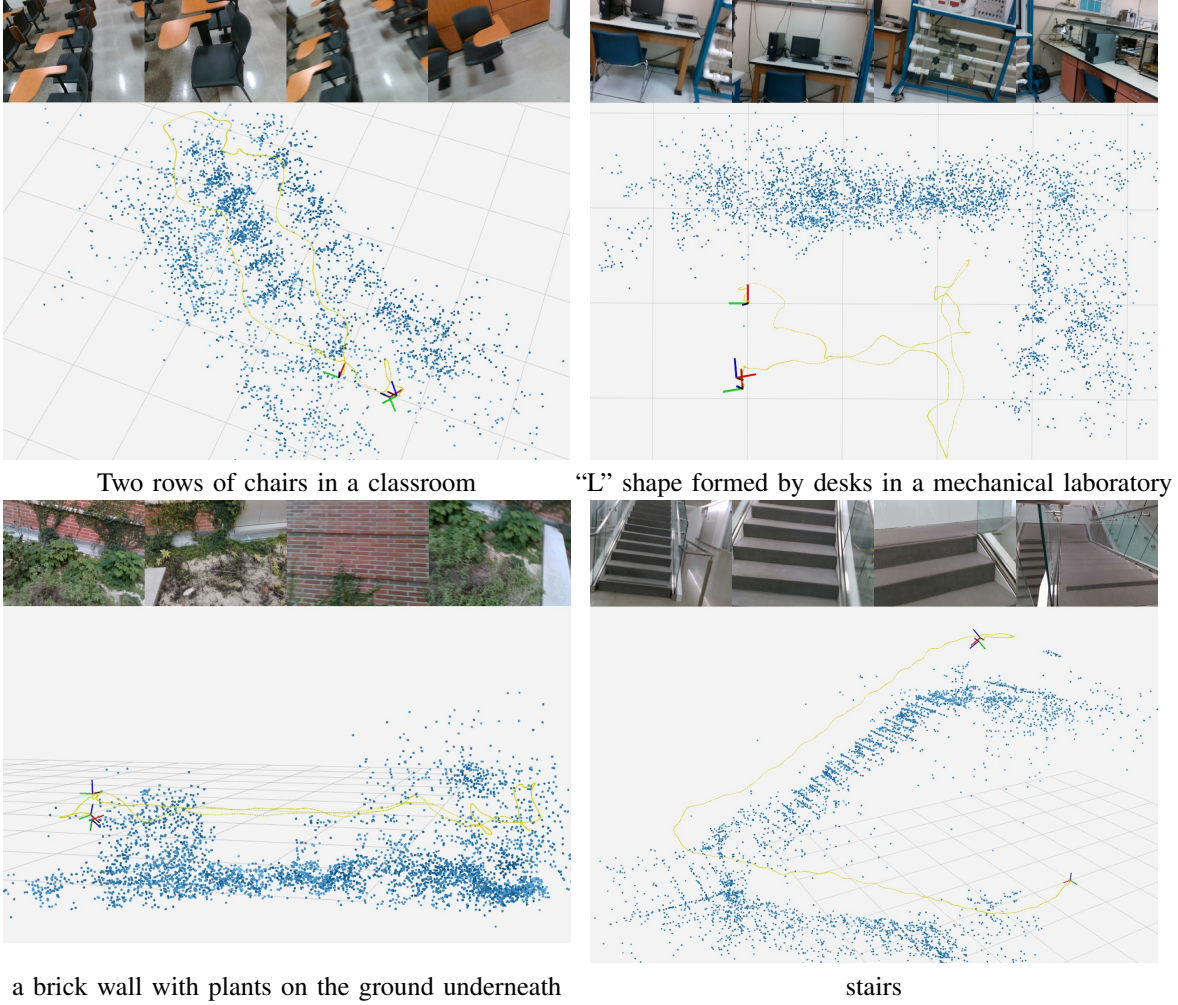


Fig. 9. Sample sequences in VOID dataset (best viewed in color at $5\times$). In each panel, the top inset shows 4 sample images of a video sequence in our VOID dataset; the bottom shows the sparse pointcloud reconstruction (blue) and camera trajectory (yellow) from our VIO.

APPENDIX B MORE RESULTS ON VOID DATASET

In the main paper, we evaluated our approach on the VOID depth completion benchmark in Sec. VII-C, and Sec. VII-D provided quantitative results in Table V and IV and qualitative results in Fig. 7. Here, we provide additional qualitative results in Fig. 10 to show how our approach performs on a variety of scenes – both indoor and outdoor – from the VOID dataset. The figure is arranged in two panels of 3×2 grids, where each panel contains a sample RGB image (left) that is fed to our depth completion network as input, and the corresponding colored pointcloud (right) produced by our approach, viewed at a different vantage point. The pointclouds are obtained by back-projecting the color pixels to the estimated depth. We used an input sparse depth density level of $\approx 0.5\%$ to produce the results. Our approach can provide detailed reconstructions of scenes from both indoor (e.g. right panel, last row: equipment from mechanical lab) and outdoor settings (e.g. left panel: flowers and leaves of plants in garden). It is also able to recover small objects such as the mouse on the desk in the mechanical lab, and structures at very close range (e.g. left panel, last row: staircase located less than half a meter from the camera).



Fig. 10. *Qualitative results on VOID dataset.* In each panel, the left shows a sample RGB image fed to our depth completion network as input; the right shows the completed depth map back-projected to 3D, colored, and viewed from a different vantage point. Our method recovers the scene structure with details at various ranges in both indoor and outdoor settings.

TABLE VI
QUANTITATIVE POSE ABLATION STUDY KITTI ODOMETRY SEQUENCE 09 AND 10.

Pose	ATE (<i>m</i>)	ATE-5F (<i>m</i>)	RPE (<i>m</i>)	RRE ($^{\circ}$)
<i>Sequence 09</i>				
Euler	34.38	0.091	0.107	0.176
Exp.	27.57	0.091	0.108	0.170
Exp. w/ Consistency	18.18	0.080	0.094	0.157
<i>Sequence 10</i>				
Euler	32.37	0.067	0.094	0.251
Exp.	25.18	0.059	0.091	0.225
Exp. w/ Consistency	24.60	0.059	0.081	0.218

We perform an ablation study on our pose representation by jointly training our depth completion network and pose network on KITTI depth completion dataset and testing only the pose network on KITTI Odometry sequence 09 and 10. We evaluate the performance of each pose network using metrics described in Sec. C-A. While performance of exponential parameterization and Euler angles are similar on ATE-5F, and RPE, exponential outperforms Euler angles in ATE and RRE on both sequences. Our model using exponential with pose consistency performs the best.

APPENDIX C POSE ABLATION STUDY

In the main paper, we focus on the depth completion task and hence we evaluate the effects of different pose parameterizations and our pose consistency term by computing error metrics relevant to the recovery of the 3D scene on both the VOID and KITTI depth completion benchmarks. Here, we focus specifically on pose by directly evaluating the pose network on the KITTI odometry dataset in Table VI. We show qualitative results on the trajectory obtained by chaining pairwise camera poses estimated by each pose network in Fig. 11 and provide an analysis of the results in Sec. C-B.

A. Pose Evaluation Metrics

To evaluate the performance of the pose network and its variants, we adopt two most widely used metrics in evaluating simultaneous localization and mapping (SLAM) systems: absolute trajectory error (ATE) and relative pose error (RPE) [31] along with two novel metrics tailored to the evaluation of pose networks.

Given a list of estimated camera poses $\hat{g}^T \doteq \{\hat{g}_1, \hat{g}_2, \dots, \hat{g}_T\}$, where $\hat{g}_t \in SE(3)$, relative to a fixed world frame, and the list of corresponding ground truth poses $g^T \doteq \{g_1, g_2, \dots, g_T\}$, where $g_t \in SE(3)$, ATE reads

$$\text{ATE}(\hat{g}^T, g^T) = \sqrt{\frac{1}{T} \sum_{t=1}^T \|\text{trans}(g_t^{-1} \hat{g}_t)\|_2^2} \quad (8)$$

where the function $\text{trans} : SE(3) \mapsto \mathbb{R}^3$ extracts the translational part of a rigid body transformation. ATE is essentially the root mean square error (RMSE) of the translational part of the estimated pose over all time indices. [18] proposed a “5-frame” version of ATE (ATE-5F) – the root mean square of ATE of a 5-frame sliding window over all time indices, which we also incorporate.

While ATE measures the overall estimation accuracy of the whole trajectory – suitable for evaluating full-fledged SLAM systems where a loop closure module presents, it does not faithfully reflect the accuracy of our pose network since 1) our pose network is designed to estimate pairwise poses, and 2) thus by simply chaining the pose estimates overtime, the pose errors at earlier time instants are more pronounced. Therefore, we also adopt RPE to measure the estimation accuracy locally:

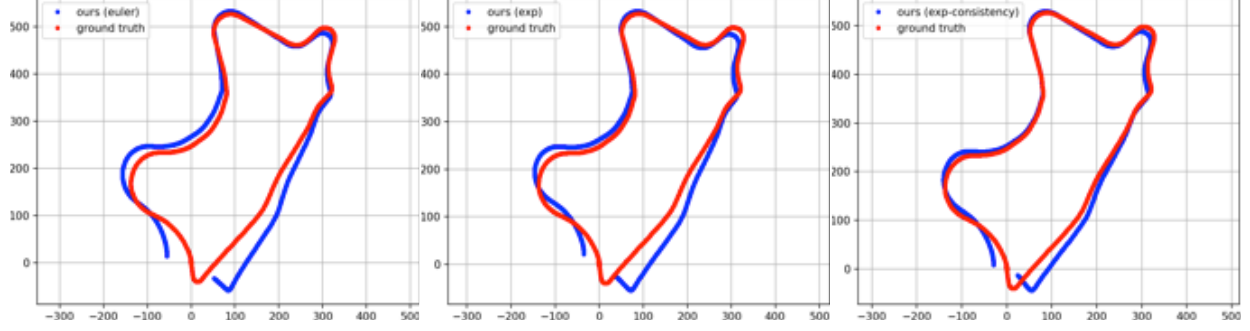
$$\text{RPE}(\hat{g}^T, g^T; \Delta) = \sqrt{\frac{1}{T-\Delta} \sum_{t=1}^{T-\Delta} \|\text{trans}((g_t^{-1} g_{t+\Delta})^{-1} (\hat{g}_t^{-1} \hat{g}_{t+\Delta}))\|_2^2} \quad (9)$$

which is essentially the end-point relative pose error of a sliding window averaged over time. By measuring the end-point relative pose $\hat{g}_{t\tau} \doteq \hat{g}_t^{-1} \hat{g}_{t+\Delta}$, where $\tau \doteq t + \Delta$, over a sliding window $[t, t + \Delta]$, we are able to focus more on the relative pose estimator (the pose network) itself rather than the overall localization accuracy. In our evaluation, we choose a sliding window of size 1, i.e., $\Delta = 1$. However, RPE is affected only by the accuracy of the translational part of the estimated pose, as we expand the relative pose error:

$$g_{t\tau}^{-1} \hat{g}_{t\tau} = (R_{t\tau}, T_{t\tau})^{-1} \cdot (\hat{R}_{t\tau}, \hat{T}_{t\tau}) \quad (10)$$

$$= (R_{t\tau}^{\top} \hat{R}_{t\tau}, -R_{t\tau}^{\top} \hat{T}_{t\tau} + T_{t\tau}) \quad (11)$$

KITTI Odometry Sequence 09



KITTI Odometry Sequence 10

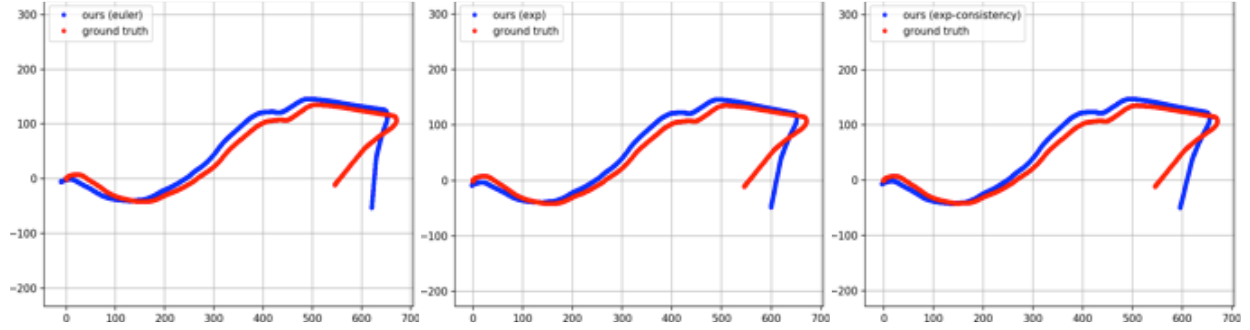


Fig. 11. *Qualitative Pose Ablation Study KITTI Odometry Sequence 09 and 10.* We perform an ablation study on our pose representation by jointly training our depth completion network and pose network on KITTI depth completion dataset and testing only the pose network on KITTI Odometry sequence 09 and 10. We obtain the camera trajectories by chaining the pairwise camera poses estimated by our pose network. We observe that the trajectory of our method using exponential parameterization trained with pose consistency (Eqn. 6) is most closely aligned with the ground-truth trajectory.

leading to $\text{trans}(g_{t\tau}^{-1}\hat{g}_{t\tau}) = -R_{t\tau}\hat{T}_{t\tau} + T_{t\tau}$, where the rotational part $\hat{R}_{t\tau}$ of the estimated pose disappears! Therefore, to better evaluate the rotation estimation, and, more importantly, to study the effect of different rotation parameterization and the pose consistency term, we propose the relative rotation error (RRE) metric:

$$\text{RRE}(\hat{g}^T, g^T; \Delta) = \sqrt{\frac{1}{T-\Delta} \sum_{t=1}^{T-\Delta} \|\log(\text{rot}(g_{t\tau}^{-1}\hat{g}_{t\tau}))\|_2^2} \quad (12)$$

where $\text{rot} : SE(3) \mapsto SO(3)$ extracts the rotational part of a rigid body transformation, and $\log : SO(3) \mapsto \mathbb{R}^3$ is the logarithmic map for rotations.

B. Ablation Study on KITTI Odometry

We perform an ablation study on the effects of our pose parameterizations and our pose consistency in Table VI and provide qualitative results showing the trajectory predicted by our pose network in Fig. 11. We jointly trained our depth completion network and our pose network on the KITTI depth completion dataset and evaluate the pose network on sequence 09 and 10 of the KITTI Odometry dataset.

For sequence 09, our pose network using exponential parameterization performs comparably to Euler angles on the ATE-5F and RPE metrics while outperforming Euler by $\approx 20\%$ on ATE and $\approx 3.4\%$ on RRE. This result suggests that while within a small window Euler and exponential perform comparably on translation, exponential is a better pose parameterization and globally more correct. We additionally see that exponential outperforms Euler angles on all metrics in sequence 10.

Our best results are achieved using exponential parameterization with our pose consistency term (Eqn. 6): on sequence 09, it outperformed Euler and exponential without pose consistency by $\approx 47.1\%$ and $\approx 28.9\%$ on ATE, $\approx 12.1\%$ and $\approx 13\%$ on RPE, $\approx 10.8\%$ and $\approx 7.6\%$ on RRE, respectively, and both by $\approx 12.1\%$ on ATE-5F. On sequence 10, it outperformed Euler and exponential by $\approx 24\%$ and $\approx 2.3\%$ on ATE, $\approx 13.8\%$ and $\approx 11\%$ on RPE, and $\approx 13.1\%$ and $\approx 3.1\%$ on RRE, respectively. It also beat Euler by $\approx 12\%$ on RPE and is comparable to exponential on the metric.

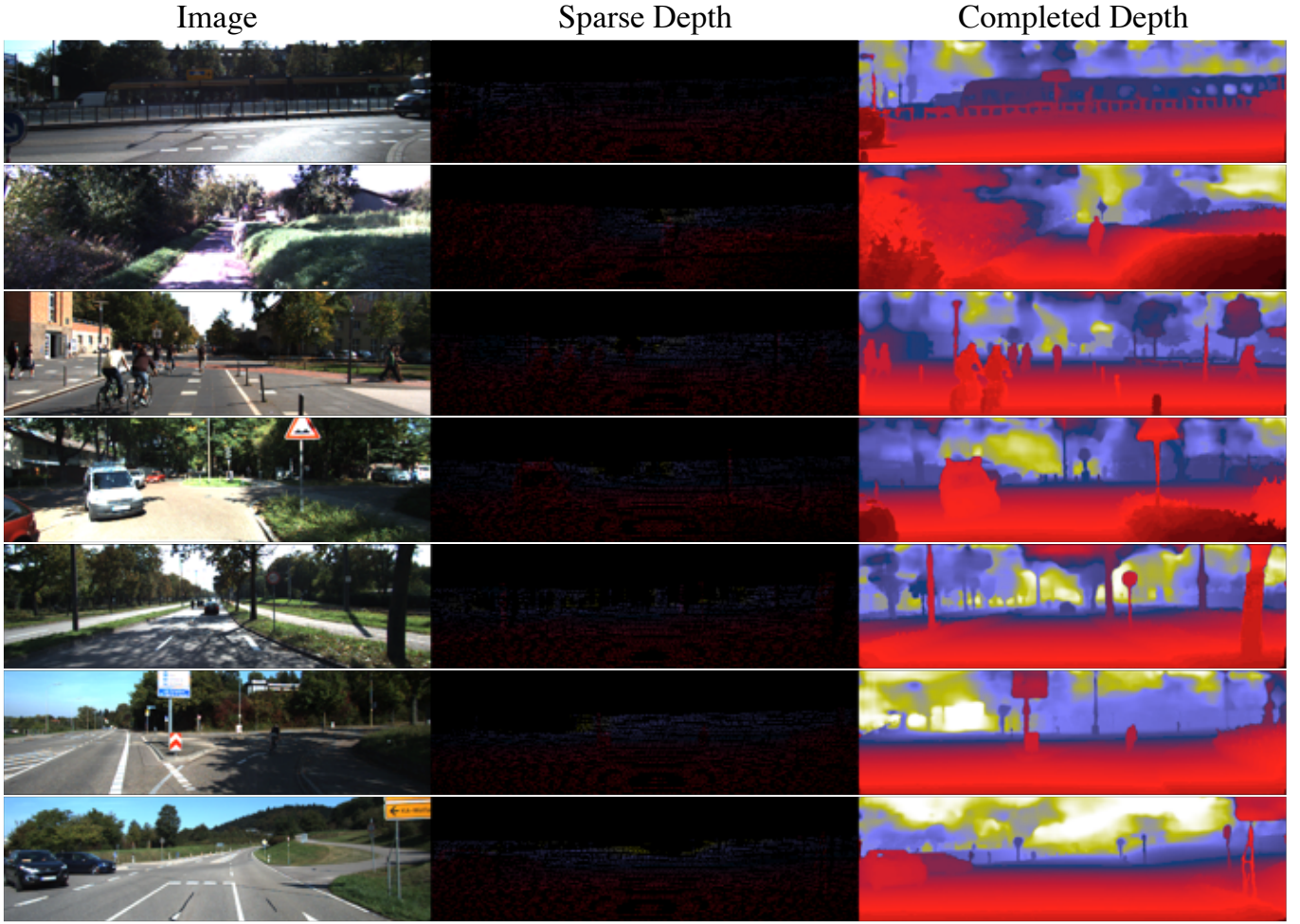


Fig. 12. *Qualitative Results on KITTI Depth Completion Test Set.* We show results from various scenes on the KITTI test set. The sparse depth input on the KITTI benchmark is concentrated on the lower half of the image domain. Our network learns to predict structures that do not have any sparse points (e.g. street sign in row 3, 5, and 6). Also, we are able to recover pedestrians (e.g. rows 2, and 3) and thin structures well (e.g. guard rails in row 1, poles in row 3, 4, 5, and 6, and 7).

APPENDIX D MORE RESULTS ON KITTI DEPTH COMPLETION BENCHMARK

In the main paper, we evaluated our approach on the KITTI depth completion benchmark test set in Sec. VII-A and performed an ablation study on the validation set in Sec. VII-B. Quantitative results are shown in Table II, III and qualitative results in Fig. 4. However, as the KITTI online depth completion benchmark only shows the first 20 samples from the test set, we provide additional qualitative results on a variety of scenes in Fig. 12 to better represent our performance on the test set.

The results in Fig. 12 were produced by our VGG11 model trained using the full loss function (Eqn. 2) with exponential parameterization for rotation. Our method is able to recover pedestrians and thin structures well (e.g. the guard rails, and street poles). Additionally, our network is also able to recover structures that do not have any associated sparse lidar points (e.g. structures located on the upper half of the image domain). This can be attributed to our photometric data-fidelity term (Sec. IV-A). As show in Fig. 3, our network first learns to copy the input scaffolding and to output it as the prediction. It later learns to fuse information from the input image to produce a prediction that includes elements from the scene that is missing from the scaffolding.

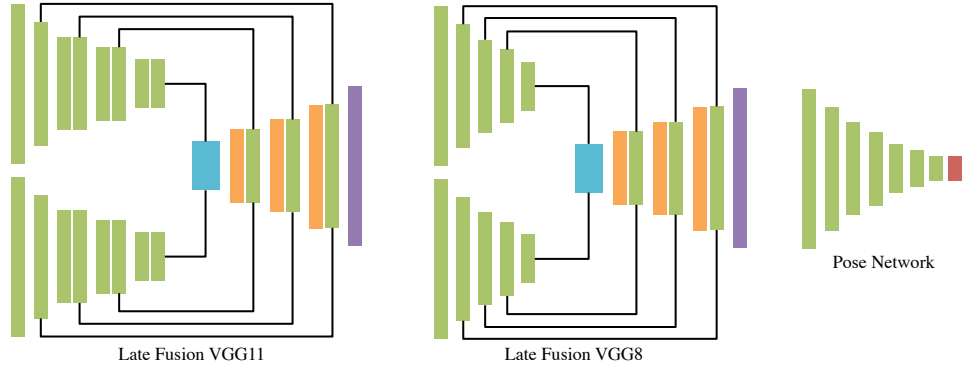


Fig. 13. Network architectures. Green denotes convolution, orange deconvolution, and purple upsampling. Blue denotes the latent representation, and red the output of pose network. Our VGG11 and VGG8 architectures following the late fusion paradigm [1], [3], and our auxiliary pose network to predict relative pose between two frames for constructing our photometric and pose consistency loss (Eqn. 4, 6). Our auxiliary pose network is used only in training and not inference.

APPENDIX E NETWORK ARCHITECTURE

We trained our model using two network architectures (Fig. 13) following the late fusion paradigm: (i) our main model using a VGG11 [32] encoder (Table VII), and (ii) our light weight model using a VGG8 [32] encoder (Table VIII). Both encoders use the same decoder (Table IX).

TABLE VII
VGG11 ENCODER ARCHITECTURE

VGG11 Encoder layer	kernel		channels		resolution		# params	input
	size	stride	in	out	in	out		
Image Branch								
conv1_image	5	2	3	48	1	1/2	≈ 3.6K	image
conv2_image	3	2	48	96	1/2	1/4	≈ 41K	conv1_image
conv3_image	3	1	96	192	1/4	1/4	≈ 166K	conv2_image
conv3b_image	3	1	192	192	1/4	1/4	≈ 331K	conv3_image
conv4_image	3	1	192	384	1/8	1/8	≈ 663K	conv3b_image
conv4b_image	3	1	384	384	1/8	1/8	≈ 1.3M	conv4_image
conv5_image	3	1	384	384	1/16	1/16	≈ 1.3M	conv4b_image
conv5b_image	3	2	384	384	1/16	1/32	≈ 1.3M	conv5_image
Depth Branch								
conv1_depth	5	2	2	16	1	1/2	≈ 0.8K	depth
conv2_depth	3	2	16	32	1/2	1/4	≈ 4.6K	conv1_depth
conv3_depth	3	1	32	64	1/4	1/4	≈ 18K	conv2_depth
conv3b_depth	3	1	64	64	1/4	1/4	≈ 37K	conv3_depth
conv4_depth	3	1	64	128	1/8	1/8	≈ 74K	conv3b_depth
conv4b_depth	3	1	128	128	1/8	1/8	≈ 147K	conv4_depth
conv5_depth	3	1	128	128	1/16	1/16	≈ 147K	conv4b_depth
conv5b_depth	3	2	128	128	1/16	1/32	≈ 147K	conv5_depth
Latent Encoding								
latent	-	-	384+128	512	1/32	1/32	0	conv5b_image conv5b_depth

Total Parameters $\approx 5.7M$

Our VGG11 [32] encoder following the late fusion paradigm [9], [3] contains $\approx 5.7M$ parameters as opposed to the $\approx 23.8M$ and $\approx 14.8M$ parameters used by [1] and [3], respectively. The || symbol denotes concatenation. Resolution ratio with respect to image size.

TABLE VIII
VGG8 ENCODER ARCHITECTURE

VGG8 Encoder layer	kernel		channels		resolution		# params	input
	size	stride	in	out	in	out		
Image Branch								
conv1_image	5	2	3	48	1	1/2	≈ 3.6K	image
conv2_image	3	2	48	96	1/2	1/4	≈ 41K	conv1_image
conv3b_image	3	2	96	192	1/4	1/8	≈ 166K	conv2_image
conv4b_image	3	2	192	384	1/8	1/16	≈ 663K	conv3b_image
conv5b_image	3	2	384	384	1/16	1/32	≈ 1.3M	conv4b_image
Depth Branch								
conv1_depth	5	2	2	16	1	1/2	≈ 0.8K	depth
conv2_depth	3	2	16	32	1/2	1/4	≈ 4.6K	conv1_depth
conv3b_depth	3	1	32	64	1/4	1/4	≈ 18K	conv2_depth
conv4b_depth	3	1	64	128	1/8	1/16	≈ 74K	conv3b_depth
conv5b_depth	3	2	128	128	1/16	1/32	≈ 147K	conv4b_depth
Latent Encoding								
latent	-	-	384+128	512	1/32	1/32	0	conv5b_image conv5b_depth
Total Parameters	≈ 2.4M							

Our light-weight VGG8 [32] encoder following the late fusion paradigm [9], [3] contains only $\approx 2.4\text{M}$ parameters as opposed to the $\approx 23.8\text{M}$ and $\approx 14.8\text{M}$ parameters used by [1] and [3], respectively. The || symbol denotes concatenation. Resolution ratio with respect to image size. Note that our light-weight model performs similarly to our VGG11 model.

TABLE IX
DECODER ARCHITECTURE

Decoder layer	kernel		channels		resolution		# params	input
	size	stride	in	out	in	out		
deconv5	3	2	512	256	1/32	1/16	$\approx 1.2\text{M}$	latent
concat5	-	-	256+384+128	768	1/16	1/16	0	deconv5 conv4b_image conv4b_depth
conv5	3	1	768	256	1/16	1/16	$\approx 1.8\text{M}$	concat5
deconv4	3	2	256	128	1/16	1/8	$\approx 295\text{K}$	conv5
concat4	-	-	128+192+64	384	1/8	1/8	0	deconv4 conv3b_image conv3b_depth
conv4	3	1	384	128	1/8	1/8	$\approx 442\text{M}$	concat4
deconv3	3	2	128	128	1/8	1/4	$\approx 147\text{K}$	conv4
concat3	-	-	128+96+32	256	1/4	1/4	0	deconv3 conv2_image conv2_depth
conv3	3	1	256	64	1/4	1/4	$\approx 147\text{K}$	concat3
deconv2	3	2	64	64	1/4	1/2	$\approx 37\text{K}$	conv3
concat2	-	-	64+48+16	128	1/2	1/2	0	deconv2 conv1_image conv1_depth
conv2	3	1	128	1	1/2	1/2	$\approx 1.2\text{K}$	concat2
output	-	-	-	-	1/2	1	0	\uparrow conv2
Total Parameters $\approx 4\text{M}$								

Our decoder contains $\approx 4\text{M}$ parameters. The || symbol denotes concatenation and the \uparrow symbol denotes upsampling. Resolution ratio with respect to image size.

TABLE X
POSE NETWORK ARCHITECTURE

Pose Network layer	kernel		channels		resolution		# params	input
	size	stride	in	out	in	out		
conv1	7	2	6	16	1	1/2	$\approx 4.7\text{K}$	image pair
conv2	5	2	16	32	1/2	1/4	$\approx 13\text{K}$	conv1
conv3	3	2	32	64	1/4	1/8	$\approx 18\text{K}$	conv2
conv4	3	2	64	128	1/8	1/16	$\approx 74\text{K}$	conv3
conv5	3	2	128	256	1/16	1/32	$\approx 295\text{K}$	conv4
conv6	3	2	256	256	1/32	1/64	$\approx 295\text{K}$	conv5
conv7	3	2	256	256	1/64	1/128	$\approx 295\text{K}$	conv6
output	3	1	256	6	1/128	1/128	$\approx 14\text{K}$	conv7
Total Parameters		$\approx 1\text{M}$						

Our auxiliary pose network contains $\approx 1\text{M}$ parameters and is only used during training to construct the photometric and pose consistency loss (Eqn. 4, 6). The output is averaged along its width and height dimensions to result in a 6 element vector – of which 3 elements are used to compose rotation and the rest for translation.

Depth completion networks. Our VGG11 and VGG8 model (Fig. 13) contain a total of $\approx 9.7\text{M}$ and $\approx 6.4\text{M}$ parameters, respectively. In comparison to [1] with $\approx 27.8\text{M}$ parameters and [3] with $\approx 18.8\text{M}$, our VGG11 model have a 65.1% and 48.4% reduction in parameters over [1] and [3], respectively; our VGG8 model have a 80% and 66% reduction over [1] and [3]. The image and depth branches of the encoder process the image and depth inputs separately – weights are not shared. The results of the encoders are concatenated as the latent representation and passed to the decoder for depth completion. The decoder makes the prediction at 1/2 resolution. The final layer of the decoder is an upsampling layer.

Pose Network. Our pose network takes a pair of images as input and regresses the relative pose between the images. Reversing the order of the image will reverse the relative pose as well. We take the average across the width and height dimensions of the pose network output to produce a 6 element vector. We use 3 elements to model rotation and the rest to model translation.

Including Pose Network in Total Parameters. We follow the network parameter computations of [3] who employs an additional network trained on ground truth for regularization during training. Our pose network (Table X) is an auxiliary network that is only used in training, and not during inference. Hence, we do not include it in the total number of parameters. However, even if we do, our pose network has $\approx 1\text{M}$ parameters, making our total for VGG11 to be $\approx 10.7\text{M}$ and VGG8 to be $\approx 7.4\text{M}$. Our VGG11 model is still has a 61.5% reduction in parameter, and our VGG8 a 73.4% over the 27.8M parameters used by [1]. If we include the auxiliary prior network of [3], containing 10.1M parameters, that is used for regularization during training, then [3] has a total of 28.8M parameters. Our VGG11 model, therefore, has a 62.8% reduction in parameters over [3] and our VGG8 has a 74.3% reduction.