# Input-level Inductive Biases for 3D Reconstruction

Wang Yifan[1*]   Carl Doersch[2]   Relja Arandjelović[2]   João Carreira[2]   Andrew Zisserman[2,3]

[1]ETH Zurich     [2]DeepMind     [3]VGG, Department of Engineering Science, University of Oxford
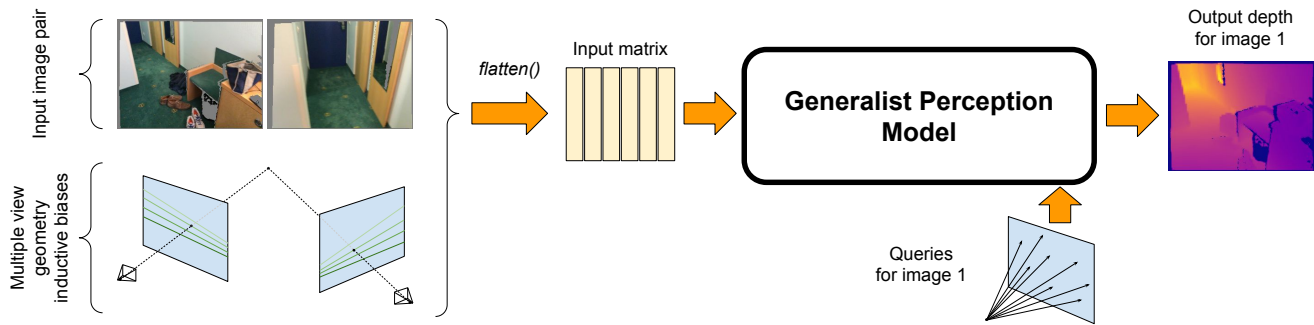
Figure 1. **Input-level inductive biases.** We explore 3D reconstruction using a generalist perception model, the recent Perceiver IO [23] which ingests a matrix of unordered and flattened inputs (*e.g.* pixels). The model is interrogated using a query matrix and generates an output for every query – in this paper the outputs are depth values for all pixels of the input image pair. We incorporate inductive biases useful for multiple view geometry into this generalist model without having to touch its architecture, by instead encoding them directly as additional inputs.

## Abstract

*Much of the recent progress in 3D vision has been driven by the development of specialized architectures that incorporate geometrical inductive biases. In this paper we tackle 3D reconstruction using a domain agnostic architecture and study how to inject the same type of inductive biases directly as extra inputs to the model. This approach makes it possible to apply existing general models, such as Perceivers, on this rich domain, without the need for architectural changes, while simultaneously maintaining data efficiency of bespoke models. In particular we study how to encode cameras, projective ray incidence and epipolar geometry as model inputs, and demonstrate competitive multi-view depth estimation performance on multiple benchmarks.*

## 1. Introduction

The focus of modern computer vision research is, to a large extent, to identify good architectures for each task of interest. There are many tasks of interest, ranging from classical ones such as optical flow [19], to highly specialized (yet arguably important) ones such as recognizing equine action units [31]. Creating dedicated models for every possible task naturally results in a sprawling catalog of architectures.

Eventually it seems desirable to build a more general visual system that can deal with most perceptual problems. To get there, one option is to combine state-of-the-art systems for all of those problems, but this would be complex, inelegant and not scalable. Another option is to employ models without much customization or inductive biases for any particular task, but these models will by definition be less data-efficient and hence less accurate than specialized ones given a fixed data budget.

In this paper we explore the single-general-model route. We ask the following question: can the lack of architecture-level inductive biases be replaced by extra inputs which encode our knowledge about the problem structure? In other words, can we feed those priors as inputs rather than hard-wire them into the model architecture (Fig. 1), like a loadable software solution instead of a more rigid hardware solution. As the general model we employ the recently published Perceiver IO [23] and as domain we focus on multi-view geometry and 3D reconstruction, an area of computer vision where architectural specialization is particularly exuberant [20, 22, 30, 36, 44, 66, 69].

Our main contribution is in mapping out and evaluating some of the options for expressing priors for 3D reconstruction as input features, in particular in the setting of depth estimation from stereo image pairs. We consider concepts in multiview geometry such as camera viewpoint, light ray direction and epipolar constraints. Similar to the prior work we compare with [20, 22, 30, 56], we assume ground truth

---

[*]Work done during internship at DeepMind.

cameras are given, but they could in principle be computed by the model as well and passed back as inputs recurrently.

We experiment on multiple datasets—ScanNet [7], SUN3D [59], RGBD-SLAM [49] and Scenes11 [56]—and present results that are comparable or better to those obtained by state-of-the-art specialized architectures on all of them. This is achieved without using cost volumes, warping layers, etc., and in fact (proudly) without introducing any architectural innovation. Instead, we propose powerful input-level 3D inductive biases that substantially improve data efficiency. This paper reflects a new avenue for problem solving in computer vision, in which domain knowledge is valued but applied in a flexible manner, as additional model inputs.

## 2. Related Work

Our work is part of a long trend in computer vision of simplifying and unifying architectures. It was noted a decade ago that big data along with simple architectures are "unreasonably effective" [15] at solving many perception problems, and subsequent progress has only reinforced this [53]. Computer vision has moved from architectures like ConvNets, which are highly general image processors [29], to methods that are based on Transformers [57] such as ViT [10] and Perceivers [23, 24], where the underlying Transformer can be equally effective across multiple domains like sound and language. Unifying architectures is useful because architectural improvements can be propagated across tasks and domains trivially. It also enables sharing and transferring information across modalities and tasks [48, 65], which is critical for tasks with little data.

However, seeking general-purpose architectures does not mean we should discard insights about geometry when solving a geometric problem. Decomposing the problem into feature matching and triangulation was an early component of stereo systems [18, 41]. More recent systems have relied on learning, especially for learning descriptors which are compared across images to find correspondence, either by directly searching for matches across images [3, 33, 35, 60] or by computing 4D correlation volumes [2, 5, 13, 27, 61, 62, 66], or a combination [14]; scaling these methods can be problematic as the number of considered matches grows. Several recent works [17, 32, 63] inferred correspondences by aggregating sample points along the epipolar line with a transformer; however matches are still represented and sampled explicitly. Similar to our work, Cam-Convs [12] leveraged input-level geometric priors (camera intrinsics) for more robust single-view depth estimation under variable camera. Our work considers a more general application - multi-view depth estimation, where we include also the camera relative poses and the epipolar embedding.

The broader field of correspondence learning has a variety of approaches for integrating global and local inference. Early approaches to deep optical flow and correspondence estimation [8, 11] used direct regression, as our approach does, but later works found correlations and cost volumes more effective [11, 50, 55]. Perceiver IO [23], however, shows strong flow performance with direct regression. Transformers have also contributed to improvements in more general scene correspondence [26, 51, 58], and even using learned correspondence to improve few-shot learning [9], though these transformers are still applied on feature grids with relatively complex mechanisms to represent correspondence explicitly. The grids are taken from prior work on correspondence that uses deep learning, where explicit pairwise comparisons and cost volumes are a staple of top-performing methods [6, 36, 39, 43–45, 64].

Our work also sits in the broader field of deep learning for 3D reconstruction, where there have been a wide variety of proposals for representing 3D inductive biases. Early works like DeepTAM [69] emphasize the importance of representing per-image depth maps and rays. More recent works have made use of deep implicit models to represent 3D [4, 37, 42], introducing the idea that deep representations should be queried with points. While this work has been extended to more complex scenes in NeRF [38] and its many derivatives, these typically require many images of the same scene and an expensive offline training process. Online methods typically rely on more explicit but expensive 3D representations like voxel grids [25, 40, 52, 67]. Particularly relevant is TransformerFusion [1], which uses Transformers to attend from its voxel grid representation to the input images, although this approach still suffers from problems with memory and resolution due to the voxel grid.

### 2.1. Review of Perceiver IO

For the general perception model we use Perceiver IO [23], and we briefly review it here. The model is based on Transformers [57] in that it treats its input as a simple series of tokens, and attention is the main workhorse. First, cross-attention is performed between the input tokens and a fixed-size set of internal vectors ('latents'), thus obtaining a compressed representation of the input. Then, a series of self-attentions is performed within the latents, enabling this architecture to scale well to large inputs (*e.g.* high resolution images) and to stack many layers without hitting memory issues, since there are much fewer latents than input tokens. The final step is another cross-attention, this time between a set of externally specified 'queries' and the latents, which produces an output array of desired size (one element for each query). Queries are typically some encoding of pixel position and are quite dense (*e.g.* one per pixel). The architecture achieves strong results on a large variety of tasks and domains, such as image classification, optical flow, natural language understanding, and StarCraft II, making it a nat-
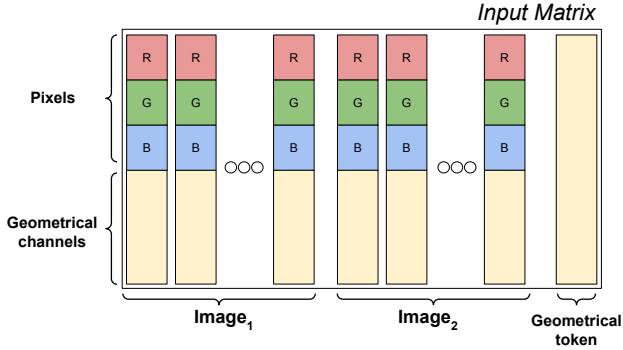
Figure 2. **Geometrical Embeddings.** The input to a Perceiver model is a matrix. Instead of vanilla positional embeddings, we introduce geometrical embeddings that encode inductive biases from multiple view geometry. We form the input matrix by concatenating pixel values with these embeddings: as extra per-pixel channels and/or as extra tokens.

ural fit for the general perception model used in this work (Fig. 1).

# 3. Featurizing Multiple View Geometry

In this section, we demonstrate how to inject geometric inductive biases into a general perception model, Perceiver IO [23] (Sec. 2.1), without changing its architecture. We consider the case of 3D reconstruction from an image pair – the inputs are pixels and calibrated cameras, and the output is depth at each pixel.

If we follow prior work, such as the optical flow network from Perceiver IO [23], then we can treat each pixel (or, more generally, each vector in a feature grid) as an input element. We then tag each pixel with an encoding for its position within the grid as input, and potentially with an additional tag to indicate which of the two input images the pixel belongs to. The output could be processed similarly: we use the same tagged pixels (or features) as queries in order to get a depth value for each pixel.

In practice, however, we expect this approach to overfit given the relatively small datasets that are available for training geometric inference. A high-capacity model can easily memorize the depth for each image, rather than learning a procedure which matches features across images and performs triangulation in a way that can generalize to unfamiliar scenes.

Our hypothesis is that we can create a more data-efficient learning algorithm by simply providing the Perceiver IO with information that describes the geometry as input. In the ideal case, Perceiver IO can learn to use this information correctly without the computational pipeline being prescribed by a complex, restrictive architecture.

In particular, we explore providing information that lets the network more easily 1) represent 3D space to allow triangulation, and 2) find correspondences, which are the two

main components of any general stereo system. Towards 1, we explore providing camera information in the form of encoded camera matrices, as well as the encodings of the rays at every pixel. Towards 2, we encode epipolar planes for each pixel, which tells the network which pixels might be in correspondence. Our main contribution in this work is to show that together, these geometric quantities can improve the inferred 3D geometry without any changes to the network architecture.

The geometric information is provided as input to the network. We explore two main ways to operationalize this (Fig. 2): 1) by fusing the information with all input elements via concatenating it along the channel dimension, and 2) by expanding the input set with additional 'geometrical' tokens.

## 3.1. Featurizing Cameras

A camera is one of the most important components for multiview geometry, providing the necessary information to perform triangulation [16]. We assume the commonly used pin-hole camera model parameterized with the intrinsic parameters $\boldsymbol{K} \in \mathbb{R}^{3 \times 3}$, which define the transformation from camera coordinates to image coordinates, and extrinsic parameters $\left[ \boldsymbol{R} \in \mathbb{R}^{3 \times 3}, \boldsymbol{t} \in \mathbb{R}^3 \right]$, defining the 6-DOF camera pose transformation from the world coordinates to the camera coordinates. In practice the intrinsic parameters can be obtained by off-the-shelf calibration methods [68] and the extrinsic parameters can be estimated using structure-from-motion algorithms such as COLMAP [47].

Next, we consider two alternatives to encode the camera parameters, the first one is based on constructing viewing rays connecting the camera with each pixel, and the second is directly providing the projection matrix that maps the 3D world coordinates to the 2D pixel coordinates.

**Option 1: Rays and camera center.** Let $\boldsymbol{x}_{j,i} \in \mathbb{R}^2$ be the image coordinate of pixel $i$ in image $j$. It can be uniquely represented in the 3D space using the viewing ray, which can be further parameterized using the camera center, $\boldsymbol{c}_j \in \mathbb{R}^3$, and the unit-length ray direction, $\boldsymbol{r}_{j,i} \in \mathbb{R}^3$ (Fig. 3). The projection matrix for the camera $j$ is a $3 \times 4$ matrix $\boldsymbol{P}_j = \boldsymbol{K}_j \left[ \boldsymbol{R}_j | \boldsymbol{t}_j \right]$. In homogeneous coordinates, the camera center $\tilde{\boldsymbol{c}}_j = [\boldsymbol{c}_j, 1]^\top$, satisfies $\boldsymbol{P}_j \tilde{\boldsymbol{c}}_j = \boldsymbol{0}$. Writing the projection matrix as $\boldsymbol{P}_j = [\boldsymbol{K}_j \boldsymbol{R}_j | \boldsymbol{K}_j \boldsymbol{t}_j]$, the camera center in the world coordinate system is

$$\boldsymbol{c}_j = -(\boldsymbol{K}_j \boldsymbol{R}_j)^{-1} \boldsymbol{K}_j \boldsymbol{t}_j = -\boldsymbol{R}_j^{-1} \boldsymbol{t}_j. \quad (1)$$

The unnormalized viewing ray direction can be computed as

$$\bar{\boldsymbol{r}}_{j,i} = (\boldsymbol{K}_j \boldsymbol{R}_j)^{-1} \begin{bmatrix} \boldsymbol{x}_{j,i} \\ 1 \end{bmatrix}, \quad (2)$$

since $\boldsymbol{P}_j \left[ \bar{\boldsymbol{r}}_{j,i}, 0 \right]^\top = [\boldsymbol{x}_{j,i}, 1]^\top$, which we normalize to unit length to obtain $\boldsymbol{r}_{j,i}$.

Instead of providing $c_j$ and $r_{j,i}$ to the network in their raw form as 3-D vectors, we embed them to higher dimensional Fourier features, as it was shown empirically that this higher-dimensional encoding is better suited for further processing by neural networks [38, 57]. This is done by applying element-wise mapping $x \mapsto [x, \sin(f_1 \pi x), \cos(f_1 \pi x), \cdots, \sin(f_K \pi x), \cos(f_K \pi x)]$, where $K$ is the number of Fourier bands, and $f_k$ is equally spaced between 1 and $\frac{\mu}{2}$, with $\mu$ corresponding to the sampling rate. The sampling rate $\mu$ and number of bands $K$ are hyperparameters which can be set separately for $c_j$ and $r_{j,i}$. As a result, we obtain $6K_c + 3$ and $6K_r + 3$ Fourier features for $c_j$ and $r_{j,i}$ respectively.

**Option 2: Pixel coordinates and projection matrix.** Alternatively, since the 3D position of each pixel can be determined up to an unknown depth solely using the projection matrix $P_j$, we can also uniquely embed each pixel directly with $P_j$ and the pixel coordinate $x_{j,i}$. To this end, we flatten $P_j$ to a 12-dimensional vector, then we map this 12-D vector as well as the 2-D $x_{j,i}$ again to Fourier features using $K_{\text{matrix}}$ and $K_x$ bands, and $\mu_{\text{matrix}}$ and $\mu_x$ sampling rates, respectively. The resulting $24K_{\text{matrix}} + 12$ and $4K_x + 2$ vectors uniquely determine the geometry for a given pixel.

**Injecting the camera information into the input of Perceiver IO.** The above geometric embeddings contain all the necessary information for the network to triangulate the pixels. We now consider how to provide this information to the general perception model. Notice that in either aforementioned options, there is a camera-specific part which is identical to all pixels in the a given image, namely $c_j$ and $P_j$, and a pixel-specific part that is unique to each pixel, namely $r_{j,i}$ and $x_{j,i}$. The pixel-specific part is most naturally incorporated by concatenating it with the pixel's RGB values along the channel dimension ('geometrical channels' in Fig. 2). There are two ways of assembling the camera-specific part – again as *geometrical channels*, or as additional *geometrical tokens*.

The first way consists of simply duplicating the camera-specific embedding for all the pixels of the corresponding image and again concatenating it along the channel dimension as *geometrical channels*. This results in a total of $2 \times H \times W$ inputs of $(D_{\text{rgb}} + D_{\text{pix}} + D_{\text{cam}})$ dimensions, where $(H, W)$ is the image dimension and $D_{\text{rgb}}, D_{\text{pix}}$ and $D_{\text{cam}}$ are the total dimensions of the RGB-based inputs, pixel-specific and camera-specific geometry embeddings respectively.

Alternatively, we can treat the camera-specific embedding as a separate input *geometric token*, alongside the per-pixel inputs, yielding a total of $2 \times (H \times W + 1)$ input tokens. In order to indicate which image is a pixel associated with, we append an additional image indicator embedding to the per-pixel tokens, that is unique per image and shared among all the pixels in the same image. In our experi-
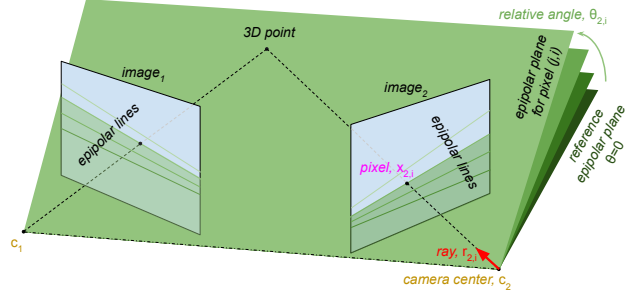


Figure 3. Geometric entities used to compute the geometrical embeddings that are passed as inputs to the perception module. For clarity, only entities related to one of the images are labeled.

ment, we encode this image indicator as a $D_{\text{ind}}$-dimensional vector using either a Fourier mapping of the image index (0/1), or as a learnable parameter. The per-image inputs contain the camera-specific geometry embedding, whereas each per-pixel input is a concatenation along the channel dimension of the RGB-based inputs, the pixel-specific geometry embedding, and the image indicator embedding. As a result, the inputs are comprised of 2 per-camera tokens with $D_{\text{cam}}$ dimensions, and $2 \times (H \times W)$ per-pixel tokens of $(D_{\text{rgb}} + D_{\text{pix}} + D_{\text{ind}})$ dimensions. Finally, to ensure the two modes of input have the same channel dimension, we pad the smaller inputs with a learnable parameter.

## 3.2. Featurizing Epipolar Cues

While the previous section uses on the geometric embedding to facilitate view triangulation, now we take one step further and exploit the given camera information to assist the search of correspondences between different images.

Correspondence estimation is paramount to multi-view geometry. The epipolar constraint is a fundamental constraint in stereo vision, which prescribes that a pair of corresponding points in the two images (projections of the 3D point) must lie on the corresponding epipolar lines, which are defined as the intersections between the image planes and the plane defined by the two camera centers and the 3D point (Fig. 3). In other words, a point in image 1 lying on epipolar line $l_1$ can only match to a point in image 2 that lies on the corresponding epipolar line $l_2$. Therefore, for a known camera pair one can compute the corresponding epipolar lines which can be used to restrict the search for point correspondences. This enables a drastically faster search while reducing the possibility of having outliers.

Similarly to camera information, the epipolar constraint is typically applied explicitly, *e.g.* by restricting the correspondence search only along epipolar lines. Instead, we provide the epipolar constraint directly as a part of the network inputs by tagging each pixel with its epipolar plane. Note that each pixel is assigned to a single epipolar plane, apart from a special case which occurs only when the projection of the other camera (the epipole) falls inside the im-

4

age, since all epipolar planes pass through the epipole; however, this degeneracy only potentially appears at a single pixel, and it is practically impossible for the epipole to align exactly with a pixel center, making this a non-issue. Next, we consider two parameterizations of the epipolar plane.

The first option encodes the *normal vector* of the epipolar plane, which can be easily computed as the normalized cross-product of $c_2 - c_1$ and $r_{j,i}$, where $r_{j,i}$ is the ray direction in (2). Formally, for pixel $i$ in image $j$, the normal vector, $n_{j,i}$, is:

$$v_{j,i} = (c_2 - c_1) \times r_{j,i} \tag{3}$$

$$n_{j,i} = \text{sign}\left([v_{j,i}]_x\right) \frac{v_{j,i}}{\|v_{j,i}\| + \epsilon}, \tag{4}$$

where the $[v_{j,i}]_x$ is the $x$-coordinate of $v_{j,i}$ and the sign disambiguates the direction of the normal vectors (opposite normals denote the same plane).

The second option parameterizes the epipolar plane as a *relative angle*, $\theta_{j,i}$, between the epipolar plane and an arbitrarily chosen reference epipolar plane, where the angles are scaled such that $\theta_{j,i} \in [-1, 1]$ (Fig. 3):

$$\theta_{j,i} = 2\left(\frac{1}{\pi} \arccos\left(\frac{n_{j,i}^\top n_{\text{ref}}}{\|n_{j,i}^\top n_{\text{ref}}\|}\right) - 0.5\right). \tag{5}$$

We choose the reference epipolar plane, which is fixed for both frames, as the plane associated with a randomly chosen pixel from the first image.

Finally, for both parametrizations, the pixel-specific epipolar encodings, $n_{j,i}$ or $\theta_{j,i}$, are embedded into Fourier features and treated as 'geometrical channels' (Fig. 2), *i.e.* concatenated to the per-pixel inputs along the channel dimension.

The epipolar embedding does not add new information compared to camera geometric embeddings described in Sec. 3.1, but it provides an additional guidance for the network to more efficiently leverage correspondence.

# 4. Experiments

We evaluate our geometrical embeddings with the Perceiver IO model on the task of depth estimation from pairs of views, a central computer vision task.

**Data.** We use the ScanNet [7] and DeMoN [56] datasets for training and testing. For ScanNet we use the frame selection as provided by [30], which yields 94212 training pairs and 7517 test pairs. The DeMoN dataset combines SUN3D [59], RGBD-SLAM [49] and Scenes11 [56]. It has a total of of 166,285 training image pairs from 50420 scenes and 288 test image pairs. Both datasets contain invalid depth measurements, following the community common practice, we mask the depths out of [0.1, 10] as invalid.

**Implementation details.** We train our model with the commonly used L1LOG loss [21], $\mathcal{L}(d, d^*) = |\log(d) - \log(d^*)|$, where $d$ and $d^*$ are the predicted and ground truth depth values. Unless otherwise stated, we process images at $240 \times 320$ resolution. The raw RGB values are transformed to 64-d (*i.e.* $D_{\text{rgb}} = 64$) color features by a standard convolutional preprocessor described in Perceiver IO [23], which consists of 1-layer convolution with receptive field 7 and stride 2, followed by batch normalization, ReLU and stride-2 max-pooling, resulting in a feature grid of dimension $60 \times 80 \times 64$ for each image. Thess feature grids are combined with the geometric embeddings to form the inputs to the Perceiver IO model. We use a small version of the original Perceiver IO architecture, which uses a $2048 \times 512$ matrix for the latent representation, 1 cross-attention for the input, followed by 8 self-attention layers and 1 cross-attention for the output, where the self-attention uses 8 heads and the cross-attention has only 1. The output of the Perceiver IO model is two $60 \times 80$ depth maps. We upsample it by 4 to the original resolution with a Convex Upsampling module [55] similar to Perceiver IO applied for optical flow estimation.

For the geometric embeddings, we consider relative camera pose w.r.t. the first camera. We set $K_r = K_{\text{matrix}} = K_p = 10$, $K_o = 20$, the maximal sampling rate $\mu$ to is set to 60 for the $r_{j,i}$, $c_j$ and $P_j$, and to 120 for the epipolar cue. These hyperparameters lead to the best evaluation in our empirical study.

We apply extensive augmentations, including random color jittering, which varies the brightness, contrast saturation and hue of the RGB inputs, as well as random cropping, rotation, and horizontal flipping. We use the ADAMW [28, 34] optimizer with standard parameters $\beta_1 = 0.9$, $\beta_2 = 0.999$, and a cosine learning rate schedule without warmup, a weight decay of $1e{-}5$, a maximal learning rate of $2e{-}4$, and train for 250 epochs with a batch size of 64.

## 4.1. Geometrical Embeddings

We present results in a top-down matter, starting with the higher-level questions: are camera and epipolar geometrical embeddings useful? Are they complementary? We then trickle down and study finer-grained design decisions for each of these two families of geometrical embeddings. In this subsection all experiments are done on ScanNet. For statistical robustness, we train three models using different random seeds and report the median result.

**Coarse-grained analysis.** We consider the best-performing options (according to the fine-grained analysis in the next subsection) for camera and epipolar embeddings Secs. 3.1 and 3.2.

Tab. 1 shows that, compared to using just standard pixel positional embeddings, any of the geometric embeddings contribute to substantial depth estimation accuracy im-

| camera | epipolar | training data proportion | | |
|---|---|---|---|---|
| embedding | cue | 30% | 50% | 100% |
| | | 0.2568 | 0.2423 | 0.2340 |
| ✓ | | **0.1350** | **0.1293** | 0.1234 |
| | ✓ | 0.2084 | 0.2018 | 0.1853 |
| ✓ | ✓ | 0.1371 | 0.1304 | **0.1204** |

Table 1. The effect of inputs on training efficiency and generalization (evaluated using absolute relative difference – lower is better), using the best option for each mode.

| camera parametrization | camera assembling | epipolar parametrization | abs. rel.diff |
|---|---|---|---|
| $\boldsymbol{c}_j, \boldsymbol{r}_{j,i}$ | channel | – | **0.1234** |
| $\boldsymbol{c}_j, \boldsymbol{r}_{j,i}$ | token | – | 0.1249 |
| $\boldsymbol{P}_j, \boldsymbol{x}_{j,i}$ | channel | – | 0.1345 |
| $\boldsymbol{P}_j, \boldsymbol{x}_{j,i}$ | token | – | 0.1805 |
| $\boldsymbol{c}_j, \boldsymbol{r}_{j,i}$ | channel | $\boldsymbol{n}_{j,i}$ | 0.1235 |
| $\boldsymbol{c}_j, \boldsymbol{r}_{j,i}$ | channel | $\theta_{j,i}$ | **0.1204** |

Table 2. Comparison between different parameterization options for camera and epipolar embeddings (using absolute relative difference).

provement, with the camera embedding reducing the absolute relative difference almost by half. Interestingly, while the epipolar embedding itself does not provide sufficient information to perform triangulation, the epipolar embedding alone (row 3) can enhance the result as it provides additional guidance for correspondence estimation.

When provided with both the camera and epipolar embeddings (row 4), our model performs similarly as when using the camera embedding alone. As the amount of training data increases however, the epipolar embedding seems to start contributing positively to the overall accuracy.

**Fine-grained analysis.** We now get down to more detailed analysis and compare the different options introduced in Secs. 3.1 and 3.2. First, we compare the two proposed camera parameterizations, namely using camera center and ray direction, $\boldsymbol{c}_j$ and $\boldsymbol{r}_{j,i}$, or directly using the projection matrix and the pixel positions, $\boldsymbol{P}_j$ and $\boldsymbol{x}_{j,i}$, as well as the two approaches for assembling this information into the input via geometrical channels or geometrical tokens. As the upper part of Tab. 2 shows, using camera center and ray direction has a consistent advantage regardless of the assembling method, likely thanks to its compactness. At the same time, we observe that concatenating the geometric embedding channel-wise to the RGB inputs compares favorably with using the geometric embedding as a separate token. This is likely due to the fact that the concatenation provides a more direct association between the geometry and the pixel-wise RGB information.

Based on the best camera configuration, we evaluate the two options for the parameterization of the epipolar cue. As the lower part of Tab. 2 shows, the angle parameterization slightly outperforms the normal parameterization, likely because the randomness in choosing the reference epipolar plane reduces the overfitting.

**Queries.** We evaluate two types of queries. As the first option, the queries take the same form as the inputs, using both the RGB features and the constructed geometric embeddings. Alternatively, we also experiment with discarding the RGB and querying using only the geometric embeddings.

We show the progression of training loss and validation

error in Fig. 5. We observe from the validation curve (right), that when including the RGB information (green) in the queries, the network initially learns slightly faster, but as the training progresses, this is outperformed by the queries that contain only the geometric-embedding. On the other hand, the training loss of the RGB-included queries remains smaller than that of the RGB-excluded queries, suggesting that the RGB information eventually leads the network to overfit by overly attending to the texture information. An example of such behavior is shown in Fig. 6.

## 4.2. Comparison with State-of-the-Art Methods

We now compare our best model to the state-of-the-art on 4 different datasets: ScanNet, Sun3D, RGBD-SLAM and Scenes11. The results are shown in Tab. 3 and indicate that using geometrical embeddings with a very generic model matches and sometimes outperforms specialized state-of-the-art models. Note that the NAS model [30] uses additional ground truth normal information as supervision and enforces consistency between normals and depth.

We also evaluate the generalization ability of our method. As we demonstrate in the Appendix A, when testing on an unseen dataset of similar domain, our model performs on par with state-of-the-art methods trained specifically for that dataset, but sees a performance drop under significant domain shift. This is somewhat expected, since unlike conventional plane-sweep methods, our model doesn't have the frames aligned externally but rather learns the alignment from input cues.

## 4.3. Camera Localization

To what extent does our algorithm understand camera geometry, as opposed to simply memorizing depth [54]? One way we can find out is by asking our algorithm to perform a useful task that it was never trained to perform. It turns out that our network can actually localize cameras given 3D geometry, which is an important component of any SLAM system.

We assume that we have a pair of images and ground truth depth maps for both. We assume camera intrinsics are

6

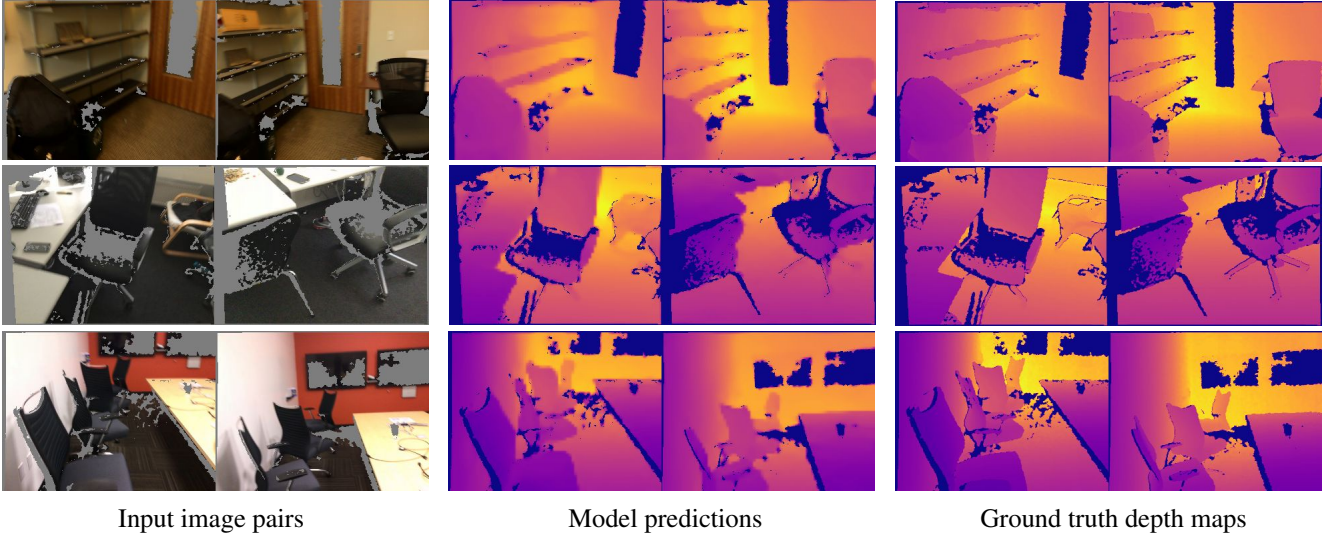| Input image pairs | Model predictions | Ground truth depth maps |

Figure 4. Examples of estimated depths using our best model on image pairs from ScanNet. Holes in the ground truth depth maps are masked out (shown in black).
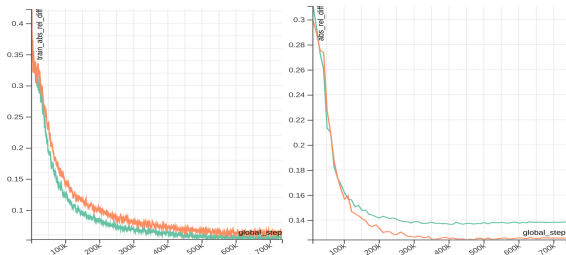


Figure 5. Queries with (green) and without (orange) RGB information. We show the training loss (left) and the validation curve (right). The change of the relative performance ranking between the two options in the validation curve suggests that overfitting to the RGB information.
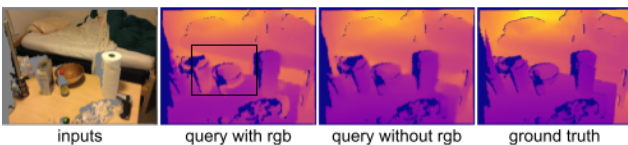


Figure 6. The effect of using RGB in the queries. Using RGB information in the query may cause artifacts due to overfitting the depth to texture.

available, but the camera position and orientation are unknown. We can randomly initialize the relative offset and orientation of the cameras, and then optimize them to minimize the L1LOG distance between the predicted depth and ground truth depth. Our underlying assumption is that errors in cameras will result in incorrect triangulation when the correspondences are correct. Therefore, if the algorithm is doing correct 3D geometry, the error should be minimized when the relative camera positions are correct. We give implementation details in Appendix B.

| dataset | methods | abs.rel. ↓ | rmse ↓ | $\delta < 1.25$ ↑ |
|---|---|---|---|---|
| ScanNet | DPSNet [22] | 0.1258 | 0.3145 | - |
| | NAS [30] | **0.1070** | **0.2807** | - |
| | IIB (ours) | 0.1159 | **0.2807** | 0.9079 |
| SUN3D | DeMoN [56] | 0.2137 | 2.4212 | 0.7332 |
| | DeepMVS [20] | 0.2816 | 0.9436 | 0.5622 |
| | DPSNet [22] | 0.1469 | 0.4489 | 0.7812 |
| | NAS [30] | 0.1271 | 0.3775 | 0.8292 |
| | IIB (ours) | **0.0985** | **0.2934** | **0.9018** |
| RGBD-SLAM | DeMoN [56] | 0.1569 | 1.7798 | 0.8011 |
| | DeepMVS [20] | 0.2938 | 0.8684 | 0.5493 |
| | DPSNet [22] | 0.1508 | 0.6952 | 0.8041 |
| | NAS [30] | 0.1314 | 0.6190 | 0.8565 |
| | IIB (ours) | **0.0951** | **0.5498** | **0.9065** |
| Scenes11 | DeMoN [56] | 0.5560 | 2.6034 | 0.4963 |
| | DeepMVS [20] | 0.2100 | 0.8909 | 0.6881 |
| | DPSNet [22] | 0.0500 | 0.4661 | 0.9614 |
| | NAS [30] | **0.0380** | **0.3710** | **0.9754** |
| | IIB (ours) | 0.0556 | 0.5229 | 0.9631 |

Table 3. Comparison with the state-of-the-art. Our method, here named IIB for *Input-level Inductive Biases*, performs competitively with these more specialized methods, doing best on two of the four datasets and coming close in the other two.

We evaluate on the SUN3D validation set. We treat the first camera as fixed at the origin, and evaluate the position of the second camera relative to it. Following prior work on camera localization [46], we evaluate two metrics. First is translation error in cm, which is simply $\|c_{est} - c_{gt}\|_2$ where $c_{est}$ and $c_{gt}$ are the estimated and ground truth camera centers, respectively. Second is rotation error in degrees, which is computed as $\arccos((\text{trace}(R_{gt}^{-1} R_{est}) - 1)/2)$, where
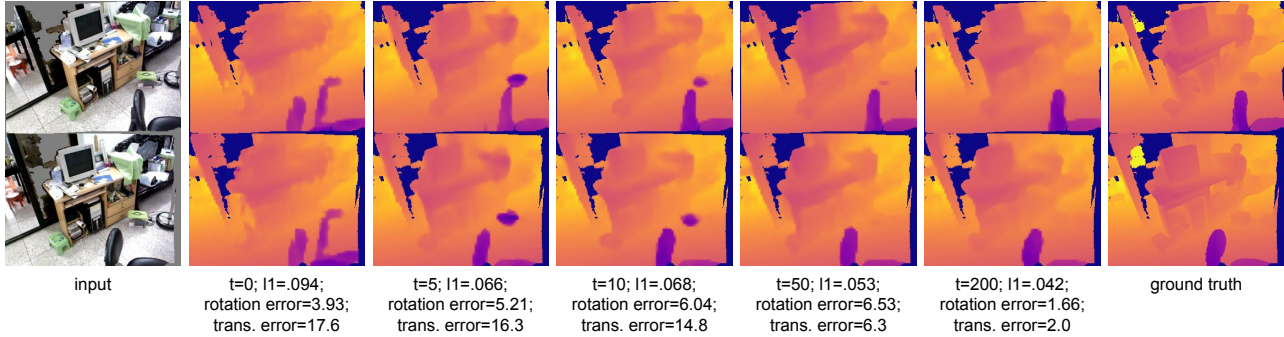
Figure 7. Progression of our iterative camera localization algorithm across 200 timesteps of optimization for the input image (left) and ground-truth depth that we are fitting to (right). At each step we show the optimization timestep (t), l1 loss between the predicted depth and the ground truth depth, as well as the rotation and translation error between the estimated and ground truth cameras (which is not used in the optimization).

| | mean rotation error (°) | median rotation error (°) | mean translation error (cm) | median translation error (cm) |
|---|---|---|---|---|
| Identity | 9.11 | 7.61 | 17.7 | 12.7 |
| Rand. init | 9.18 | 7.68 | 17.8 | 12.8 |
| Optimized | 6.67 | 4.38 | 2.5 | 1.9 |

Table 4. Camera Localization performance on SUN3D. Lower is better.

$R_{est}$ and $R_{gt}$ are the estimated and ground truth rotation matrices respectively. This is the minimum rotation angle required to align both rotations.

Tab. 4 shows our results. We report both the mean and median across examples in the dataset. We see non-trivial improvements in both metrics, localizing the cameras within a few centimeters and a few degrees of their true locations. While we don't expect this to be competitive with SOTA SLAM systems (which typically integrate information across many more images), these results clearly show that the algorithm is using geometry as expected: depth error is minimized when the cameras are at the correct locations. This may serve as a starting point for more flexible systems which don't rely on access to ground truth cameras.

Fig. 7 shows how the depth map progresses throughout 200 steps of optimizing the camera. We see that the initial depth is quite poor, with little definition on the desk, and a chair in the foreground which has been split in two, like double-vision in humans. These errors gradually resolve as the camera estimate gets better, with the algorithm able to correctly bind pixels across images using the geometry. Interestingly, the translation error improves faster than the rotation error, suggesting that the algorithm may be using the camera centers more than the ray angles in order to perform triangulation.

## 5. Discussion

The 3D nature of space is a key aspect of our reality and should be incorporated as priors into our visual models. Most models for 3D reconstruction currently incorporate 3D priors by tailoring the architectures. In this paper we investigated an alternative inspired by advances in modeling with Transformers: we featurize these priors and feed them as inputs to the model. We show that this incurs no sacrifice in performance and in fact we obtain results that are competitive with leading models on several datasets. Our exploration in the space of geometry parameterization is non-exhaustive, more indicative priors may be derived to simplify the 3D reasoning.

Having geometric priors as inputs also opens up new possibilities: given a pre-trained frozen model and ground truth depth, one can finetune the geometrical inputs in case they are unknown, *e.g.* for camera calibration or epipolar geometry estimation. Input-level inductive biases may also enable us to incorporate geometry into multimodal models, *e.g.*, those that jointly process sound, touch or text. In such a setting, the type of architecture engineering that is appropriate for vision would no longer apply, whereas input-level biases still could.

On the other hand, since our model needs to learn the 3D alignment, it expects the training and test data to have a similar distribution. Moreover, since the baseline architecture operates on a compressed latent space, how the network solves the problem is potentially less interpretable because there are no explicit correspondence establishing steps.

## 6. Acknowledgements

# References

[1] Aljaž Božič, Pablo Palafox, Justus Thies, Angela Dai, and Matthias Nießner. TransformerFusion: Monocular RGB scene reconstruction using transformers. *arXiv preprint arXiv:2107.02191*, 2021. 2

[2] Jia-Ren Chang and Yong-Sheng Chen. Pyramid stereo matching network. In *CVPR*, pages 5410–5418, 2018. 2

[3] Zhuoyuan Chen, Xun Sun, Liang Wang, Yinan Yu, and Chang Huang. A deep visual correspondence embedding model for stereo matching costs. In *ICCV*, pages 972–980, 2015. 2

[4] Zhiqin Chen and Hao Zhang. Learning implicit fields for generative shape modeling. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5939–5948, 2019. 2

[5] Xuelian Cheng, Yiran Zhong, Mehrtash Harandi, Yuchao Dai, Xiaojun Chang, Hongdong Li, Tom Drummond, and Zongyuan Ge. Hierarchical neural architecture search for deep stereo matching. *NeurIPS*, 33:22158–22169, 2020. 2

[6] Christopher Choy, Wei Dong, and Vladlen Koltun. Deep global registration. In *CVPR*, pages 2514–2523, 2020. 2

[7] Angela Dai, Angel X. Chang, Manolis Savva, Maciej Halber, Thomas Funkhouser, and Matthias Nießner. Scannet: Richly-annotated 3d reconstructions of indoor scenes. In *Proc. Computer Vision and Pattern Recognition (CVPR), IEEE*, 2017. 2, 5

[8] Daniel DeTone, Tomasz Malisiewicz, and Andrew Rabinovich. Deep image homography estimation. *arXiv preprint arXiv:1606.03798*, 2016. 2

[9] Carl Doersch, Ankush Gupta, and Andrew Zisserman. CrossTransformers: Spatially-aware few-shot transfer. In *NeurIPS*, 2020. 2

[10] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*, 2020. 2

[11] Alexey Dosovitskiy, Philipp Fischer, Eddy Ilg, Philip Hausser, Caner Hazirbas, Vladimir Golkov, Patrick Van Der Smagt, Daniel Cremers, and Thomas Brox. Flownet: Learning optical flow with convolutional networks. In *ICCV*, pages 2758–2766, 2015. 2

[12] Jose M Facil, Benjamin Ummenhofer, Huizhong Zhou, Luis Montesano, Thomas Brox, and Javier Civera. Camconvs: Camera-aware multi-scale convolutions for single-view depth. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11826–11835, 2019. 2

[13] Xiaodong Gu, Zhiwen Fan, Siyu Zhu, Zuozhuo Dai, Feitong Tan, and Ping Tan. Cascade cost volume for high-resolution multi-view stereo and stereo matching. In *CVPR*, pages 2495–2504, 2020. 2

[14] Xiaoyang Guo, Kai Yang, Wukui Yang, Xiaogang Wang, and Hongsheng Li. Group-wise correlation stereo network. In *CVPR*, pages 3273–3282, 2019. 2

[15] Alon Halevy, Peter Norvig, and Fernando Pereira. The unreasonable effectiveness of data. *IEEE Intelligent Systems*, 24(2):8–12, 2009. 2

[16] Richard I. Hartley and Andrew Zisserman. *Multiple View Geometry in Computer Vision*. Cambridge University Press, ISBN: 0521540518, second edition, 2004. 3

[17] Yihui He, Rui Yan, Katerina Fragkiadaki, and Shoou-I Yu. Epipolar transformers. In *Proceedings of the ieee/cvf conference on computer vision and pattern recognition*, pages 7779–7788, 2020. 2

[18] Heiko Hirschmuller. Stereo processing by semiglobal matching and mutual information. *IEEE TPAMI*, 30(2):328–341, 2007. 2

[19] Berthold KP Horn and Brian G Schunck. Determining optical flow. *Artificial intelligence*, 17(1-3):185–203, 1981. 1

[20] Po-Han Huang, Kevin Matzen, Johannes Kopf, Narendra Ahuja, and Jia-Bin Huang. DeepMVS: Learning multi-view stereopsis. In *CVPR*, pages 2821–2830, 2018. 1, 7

[21] Lam Huynh, Phong Nguyen-Ha, Jiri Matas, Esa Rahtu, and Janne Heikkilä. Guiding monocular depth estimation using depth-attention volume. In *European Conference on Computer Vision*, pages 581–597. Springer, 2020. 5

[22] Sunghoon Im, Hae-Gon Jeon, Stephen Lin, and In So Kweon. DPSNet: End-to-end deep plane sweep stereo. *arXiv preprint arXiv:1905.00538*, 2019. 1, 7

[23] Andrew Jaegle, Sebastian Borgeaud, Jean-Baptiste Alayrac, Carl Doersch, Catalin Ionescu, David Ding, Skanda Koppula, Daniel Zoran, Andrew Brock, Evan Shelhamer, Olivier Hénaff, Matthew M. Botvinick, Andrew Zisserman, Oriol Vinyals, and Joao Carreira. Perceiver IO: A general architecture for structured inputs & outputs. *arXiv preprint arXiv:2107.14795*, 2021. 1, 2, 3, 5

[24] Andrew Jaegle, Felix Gimeno, Andrew Brock, Andrew Zisserman, Oriol Vinyals, and Joao Carreira. Perceiver: General perception with iterative attention. In *ICML*, 2021. 2

[25] Mengqi Ji, Juergen Gall, Haitian Zheng, Yebin Liu, and Lu Fang. Surfacenet: An end-to-end 3d neural network for multiview stereopsis. In *ICCV*, pages 2307–2315, 2017. 2

[26] Wei Jiang, Eduard Trulls, Jan Hosang, Andrea Tagliasacchi, and Kwang Moo Yi. COTR: Correspondence transformer for matching across images. *arXiv preprint arXiv:2103.14167*, 2021. 2

[27] Alex Kendall, Hayk Martirosyan, Saumitro Dasgupta, Peter Henry, Ryan Kennedy, Abraham Bachrach, and Adam Bry. End-to-end learning of geometry and context for deep stereo regression. In *ICCV*, pages 66–75, 2017. 2

[28] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014. 5

[29] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. *NeurIPS*, 25:1097–1105, 2012. 2

[30] Uday Kusupati, Shuo Cheng, Rui Chen, and Hao Su. Normal assisted stereo depth estimation. In *CVPR*, pages 2189–2199, 2020. 1, 5, 6, 7

[31] Zhenghong Li, Sofia Broomé, Pia Haubro Andersen, and Hedvig Kjellström. Automated detection of equine facial action units. *arXiv preprint arXiv:2102.08983*, 2021. 1

[32] Zhaoshuo Li, Xingtong Liu, Nathan Drenkow, Andy Ding, Francis X Creighton, Russell H Taylor, and Mathias Unberath. Revisiting stereo depth estimation from a sequence-to-sequence perspective with transformers. In *ICCV*, pages 6197–6206, 2021. 2

[33] Zhengfa Liang, Yiliu Feng, Yulan Guo, Hengzhu Liu, Wei Chen, Linbo Qiao, Li Zhou, and Jianfeng Zhang. Learning for disparity estimation through feature constancy. In *CVPR*, pages 2811–2820, 2018. 2

[34] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. *arXiv preprint arXiv:1711.05101*, 2017. 5

[35] Nikolaus Mayer, Eddy Ilg, Philip Hausser, Philipp Fischer, Daniel Cremers, Alexey Dosovitskiy, and Thomas Brox. A large dataset to train convolutional networks for disparity, optical flow, and scene flow estimation. In *CVPR*, pages 4040–4048, 2016. 2

[36] Iaroslav Melekhov, Aleksei Tiulpin, Torsten Sattler, Marc Pollefeys, Esa Rahtu, and Juho Kannala. DGC-Net: Dense geometric correspondence network. In *WACV*, pages 1034–1042. IEEE, 2019. 1, 2

[37] Lars Mescheder, Michael Oechsle, Michael Niemeyer, Sebastian Nowozin, and Andreas Geiger. Occupancy networks: Learning 3D reconstruction in function space. In *CVPR*, pages 4460–4470, 2019. 2

[38] Ben Mildenhall, Pratul P Srinivasan, Matthew Tancik, Jonathan T Barron, Ravi Ramamoorthi, and Ren Ng. NeRF: Representing scenes as neural radiance fields for view synthesis. In *ECCV*, pages 405–421. Springer, 2020. 2, 4

[39] Juhong Min, Jongmin Lee, Jean Ponce, and Minsu Cho. Hyperpixel flow: Semantic correspondence with multi-layer neural features. In *ICCV*, pages 3395–3404, 2019. 2

[40] Zak Murez, Tarrence van As, James Bartolozzi, Ayan Sinha, Vijay Badrinarayanan, and Andrew Rabinovich. Atlas: End-to-end 3d scene reconstruction from posed images. In *ECCV*, pages 414–431. Springer, 2020. 2

[41] Yuichi Ohta and Takeo Kanade. Stereo by intra-and inter-scanline search using dynamic programming. *IEEE TPAMI*, pages 139–154, 1985. 2

[42] Jeong Joon Park, Peter Florence, Julian Straub, Richard Newcombe, and Steven Lovegrove. DeepSDF: Learning continuous signed distance functions for shape representation. In *CVPR*, pages 165–174, 2019. 2

[43] Ignacio Rocco, Relja Arandjelović, and Josef Sivic. Efficient neighbourhood consensus networks via submanifold sparse convolutions. In *ECCV*, pages 605–621. Springer, 2020. 2

[44] Ignacio Rocco, Mircea Cimpoi, Relja Arandjelović, Akihiko Torii, Tomas Pajdla, and Josef Sivic. Neighbourhood consensus networks. In *NeurIPS*, 2018. 1, 2

[45] Paul-Edouard Sarlin, Daniel DeTone, Tomasz Malisiewicz, and Andrew Rabinovich. SuperGlue: Learning feature matching with graph neural networks. In *CVPR*, pages 4938–4947, 2020. 2

[46] Torsten Sattler, Will Maddern, Carl Toft, Akihiko Torii, Lars Hammarstrand, Erik Stenborg, Daniel Safari, Masatoshi Okutomi, Marc Pollefeys, Josef Sivic, et al. Benchmarking 6dof outdoor visual localization in changing conditions. In *CVPR*, pages 8601–8610, 2018. 7

[47] Johannes Lutz Schönberger and Jan-Michael Frahm. Structure-from-motion revisited. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016. 3

[48] Ali Sharif Razavian, Hossein Azizpour, Josephine Sullivan, and Stefan Carlsson. Cnn features off-the-shelf: an astounding baseline for recognition. In *CVPR*, pages 806–813, 2014. 2

[49] Jürgen Sturm, Nikolas Engelhard, Felix Endres, Wolfram Burgard, and Daniel Cremers. A benchmark for the evaluation of rgb-d slam systems. In *2012 IEEE/RSJ international conference on intelligent robots and systems*, pages 573–580. IEEE, 2012. 2, 5

[50] Deqing Sun, Xiaodong Yang, Ming-Yu Liu, and Jan Kautz. Pwc-net: Cnns for optical flow using pyramid, warping, and cost volume. In *CVPR*, pages 8934–8943, 2018. 2

[51] Jiaming Sun, Zehong Shen, Yuang Wang, Hujun Bao, and Xiaowei Zhou. Loftr: Detector-free local feature matching with transformers. In *CVPR*, pages 8922–8931, 2021. 2

[52] Jiaming Sun, Yiming Xie, Linghao Chen, Xiaowei Zhou, and Hujun Bao. NeuralRecon: Real-time coherent 3D reconstruction from monocular video. In *CVPR*, pages 15598–15607, 2021. 2

[53] Richard Sutton. The bitter lesson. `http://www.incompleteideas.net/IncIdeas/BitterLesson.html`, 2019. Accessed: 2021-11-16. 2

[54] Maxim Tatarchenko, Stephan R Richter, René Ranftl, Zhuwen Li, Vladlen Koltun, and Thomas Brox. What do single-view 3d reconstruction networks learn? In *CVPR*, pages 3405–3414, 2019. 6

[55] Zachary Teed and Jia Deng. Raft: Recurrent all-pairs field transforms for optical flow. In *ECCV*, pages 402–419. Springer, 2020. 2, 5

[56] Benjamin Ummenhofer, Huizhong Zhou, Jonas Uhrig, Nikolaus Mayer, Eddy Ilg, Alexey Dosovitskiy, and Thomas Brox. Demon: Depth and motion network for learning monocular stereo. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 5038–5047, 2017. 1, 2, 5, 7

[57] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *NeurIPS*, pages 5998–6008, 2017. 2, 4

[58] Ming Wei, Ming Zhu, Yi Wu, Jiaqi Sun, Jiarong Wang, and Changji Liu. A fast stereo matching network with multi-cross attention. *Sensors*, 21(18):6016, 2021. 2

[59] Jianxiong Xiao, Andrew Owens, and Antonio Torralba. Sun3d: A database of big spaces reconstructed using sfm and object labels. In *Proceedings of the IEEE international conference on computer vision*, pages 1625–1632, 2013. 2, 5

[60] Haofei Xu and Juyong Zhang. AANet: Adaptive aggregation network for efficient stereo matching. In *CVPR*, pages 1959–1968, 2020. 2

[61] Gengshan Yang, Joshua Manela, Michael Happold, and Deva Ramanan. Hierarchical deep stereo matching on high-resolution images. In *CVPR*, pages 5515–5524, 2019. 2

[62] Jiayu Yang, Wei Mao, Jose M Alvarez, and Miaomiao Liu. Cost volume pyramid based depth inference for multi-view stereo. In *CVPR*, pages 4877–4886, 2020. 2

[63] Zhenpei Yang, Zhile Ren, Qi Shan, and Qixing Huang. Mvs2d: Efficient multi-view stereo via attention-driven 2d convolutions. *arXiv preprint arXiv:2104.13325*, 2021. 2

[64] Kwang Moo Yi, Eduard Trulls, Yuki Ono, Vincent Lepetit, Mathieu Salzmann, and Pascal Fua. Learning to find good correspondences. In *CVPR*, pages 2666–2674, 2018. 2

[65] Jason Yosinski, Jeff Clune, Yoshua Bengio, and Hod Lipson. How transferable are features in deep neural networks? *NeurIPS*, 27:3320–3328, 2014. 2

[66] Feihu Zhang, Victor Prisacariu, Ruigang Yang, and Philip HS Torr. GA-Net: Guided aggregation net for end-to-end stereo matching. In *CVPR*, pages 185–194, 2019. 1, 2

[67] Jingyang Zhang, Yao Yao, and Long Quan. Learning signed distance field for multi-view surface reconstruction. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 6525–6534, 2021. 2

[68] Zhengyou Zhang. A flexible new technique for camera calibration. *IEEE Transactions on pattern analysis and machine intelligence*, 22(11):1330–1334, 2000. 3

[69] Huizhong Zhou, Benjamin Ummenhofer, and Thomas Brox. DeepTAM: Deep tracking and mapping. In *ECCV*, pages 822–838, 2018. 1, 2

## A. Generalization on Out-Of-Domain Test Data

Our method generalizes to unseen datasets of a comparable domain. However, when testing on a significantly different domain, *e.g.* trained on indoor scenes and testing on outdoor scenes, our framework will see a performance drop. Table 5 shows the performance of a model trained on ScanNet and evaluated on RGBD and SUN3D datasets. SUN3D is similar to ScanNet, our method (trained on ScanNet only) performs reasonably well and on par with methods trained on SUN3D. RGBD datasets contain many warehouse scenes where the depth range is significantly different from the one seen during training. Our model shows overfitting in this case. This is because unlike conventional plane-sweep methods, where the network essentially computes a cost-volume from RGB features that are explicitly aligned, our method has to learn the alignment itself. How to reduce such a gap in case of domain shift is a research direction for future work. Besides introducing more augmentation techniques, we think fine-tuning and scale-normalization are some promising directions to pursue.

## B. Camera localization: implementation details

In this section, we give implementation details for Section 4.3. Given a pair of images and ground truth depth

| method | train | test | abs.rel ↓ | rmse ↓ | $\delta < 1.25$ ↑ |
|---|---|---|---|---|---|
| IIB (ours) | ScanNet | SUN3D | 0.1291 | 0.3699 | 0.8298 |
| IIB (ours) | SUN3D | SUN3D | 0.0985 | 0.2934 | 0.9018 |
| NAS [27] | SUN3D | SUN3D | 0.1271 | 0.3775 | 0.8292 |
| IIB (ours) | ScanNet | RGBD | 0.2572 | 1.3102 | 0.5101 |
| IIB (ours) | RGBD | RGBD | 0.0951 | 0.5498 | 0.9065 |
| NAS [27] | RGBD | RGBD | 0.1314 | 0.6190 | 0.8565 |

Table 5. Generalisation performance.

maps, the goal is to infer the relative position and orientation of the two cameras using the network. We use a perceiver with $c_{j,i}, r_{j,i}$ for the camera embedding and $n_{j,i}$ for the epipolar constraint.

Recall that cameras are parameterized with intrinsics $K$, and extrinsics $R$ and $t$. We assume known $K$'s, and without loss of generality that the first camera is fixed as $R_1 = I$ and $t_1 = 0$. Thus, the extrinsics of the second camera $R_2$ and $t_2$ are the optimization variables, and these are parameterized as $t_2 \in \mathbb{R}^3$ and $\hat{R}_2 \in \mathbb{R}^{3 \times 3}$. To make sure that the extrinsic matrix is a rigid transformation, we compute the final rotation matrix as $R_2 = UV$, where $U, S, V = \text{SVD}(\hat{R}_2)$, SVD is singular value decomposition, $U$ and $V$ are orthonormal, and $S$ is diagonal. Therefore we are overparameterizing $R_2$ with 9 parameters but only 3 degrees of freedom, but empirically we find this gives good results.

Our optimization objective is L1LOG loss on both images, plus regularization terms to encourage both the translation and rotation to be small. Without regularization, we find that the optimization can get stuck in local minima with very large rotation and translation, which we don't expect will make sense to the network. Specifically, we let $L_{rot}(R) = \arccos((\text{trace}(R) - 1)/2)$ in radians, and $L_{trans}(t) = \|t\|_2$. Then the final loss can be written as:

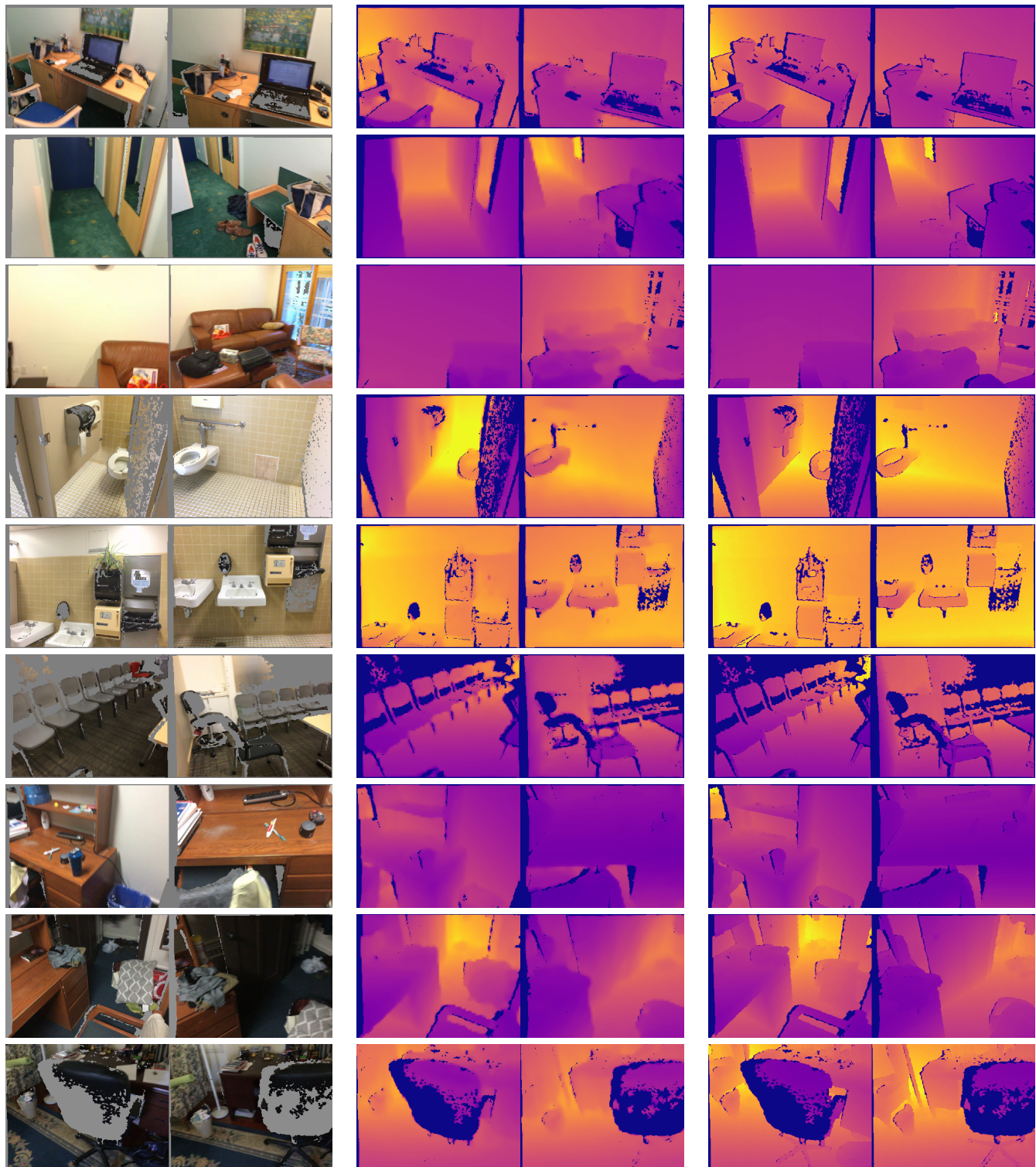$$L(R_2, t_2) = \text{L1LOG}(R_2, t_2) + \lambda_{rot} L_{rot}(R_2) + \lambda_{trans} L_{trans}(t_2)$$

We set $\lambda_{rot} = 1$ and $\lambda_{trans} = 20$. Note that during optimization, the translation coordinates use the default Sun3D scaling, which is $100\times$ smaller than what we report in our result table.

We use the ADAM optimizer with a cosine learning rate schedule for 200 steps and an initial learning rate of $5e-3$. We initialize $t_2 = \epsilon_{trans}$ where $\epsilon_{trans} \in \mathbb{R}^3$ and $\hat{R}_2 = I + \epsilon_{rot}$ where $\epsilon_{trans} \in \mathbb{R}^3$. Each element of both $\epsilon_{rot}$ and $\epsilon_{trans}$ is distributed as $\mathcal{N}(0, 0.01)$. We find that the network occasionally gets stuck in local optima, so we run with 5 different random initializations and take the solution with the best total loss $L$.

## C. Qualitative Results

Additional qualitative results are shown in Fig. 8.

11

| Input image pairs | Model predictions | Ground truth depth maps |

Figure 8. Additional examples of estimated depths using our best model on image pairs from ScanNet. Holes in the ground truth depth maps are masked out (shown in black).