# F$^3$Loc: Fusion and Filtering for Floorplan Localization

Changan Chen[1,*]    Rui Wang[2]    Christoph Vogel[2]    Marc Pollefeys[2]

[1]ETH Zürich    [2]Microsoft Mixed Reality & AI Lab Zürich

[*]Work done during his internship at Microsoft Mixed Reality & AI Lab Zürich
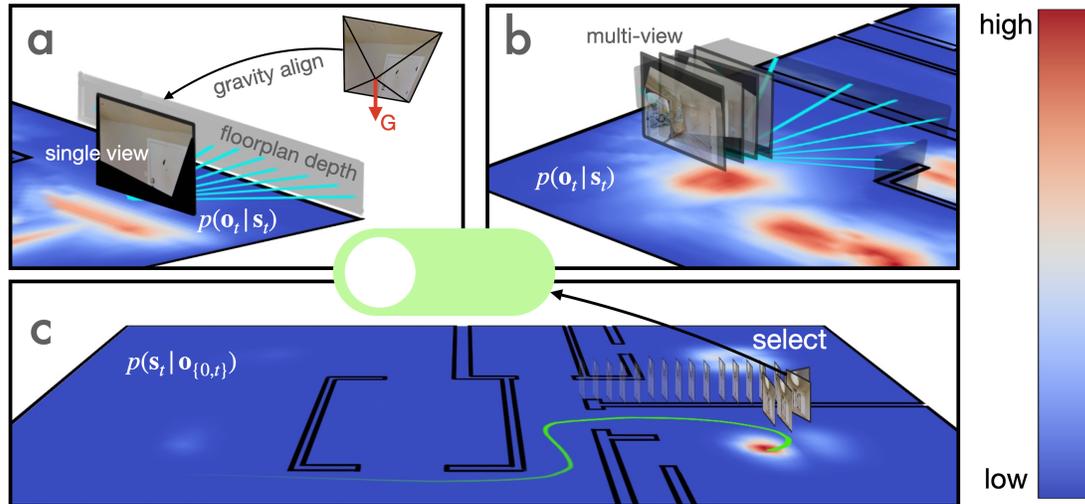
Figure 1. **Floorplan localization.** We propose a novel probabilistic model for localization within a floorplan consisting of a data-driven observation (a,b) and a temporal filtering module (c). Evidence is estimated as a 1D-range image from a single (a) and a few consecutive RGB images (b). A learned soft selection module combines the output from the complementary cues. The observation likelihood is integrated over time by an efficient SE2 histogram filter to deliver the pose posterior. Our system achieves rapid and accurate sequential localization, outperforming the state-of-the-art in recall and localization speed, while operating on consumer hardware.

## Abstract

*In this paper we propose an efficient data-driven solution to self-localization within a floorplan. Floorplan data is readily available, long-term persistent and inherently robust to changes in the visual appearance. Our method does not require retraining per map and location or demand a large database of images of the area of interest. We propose a novel probabilistic model consisting of an observation and a novel temporal filtering module. Operating internally with an efficient ray-based representation, the observation module consists of a single and a multiview module to predict horizontal depth from images and fuses their results to benefit from advantages offered by either methodology. Our method operates on conventional consumer hardware and overcomes a common limitation of competing methods [16, 17, 20, 28] that often demand upright images. Our full system meets real-time requirements, while outperform-ing the state-of-the-art [20, 28] by a significant margin.*

## 1. Introduction

Camera localization is an essential research topic in computer vision. It is key to many AR/VR applications for head-mounted or handheld mobile devices and is of great practical interest to the robotics community. Most existing works localize the camera using a pre-collected database [40][1][2] or within a pre-built 3D model [24, 34, 37–39]. However, these representations of the environment are costly in terms of storage and maintenance. In contrast, indoor environments including most commercial real estate such as warehouses, offices and apartments already possess a floorplan. The floorplan is a generic representation for indoor environments that is easily accessible, lightweight, and preserves long-term scene structure independent of a changing visual appearance, such as the furnishing of the scene. It

1

encodes rich enough information that humans can localize in an unvisited scene with its help. Therefore, we propose to localize the camera with respect to a given floorplan. This cannot only be used for indoor AR/VR applications such as floorplan navigation but also empowers robot autonomy in indoor exploration, navigation as well as search and rescue [11]. Our framework can be used complementary to indoor SLAM, where it can provide an initial guess for camera relocalization and significantly simplify detecting and verifying loop closures.

Due to its simple and compact form, floorplans contain many repetitive structures such as corners and walls. This causes ambiguity in the localization [16, 17, 28], which can be eliminated to a certain extent by using image sequences [35, 36]. However, incorporating the single frame localization into a sequential filtering framework [45, 49] is challenging. The single frame localization needs to be accurate and its efficiency is crucial to ensure a high frequency of the filter with a large amount of samples [20, 28]. To tackle these challenges, we propose a data-driven multi-view geometry based localization framework, that is both fast and accurate. Furthermore, we integrate this framework into a novel and highly efficient histogram filter that outputs a probability over poses and, thus, allows for multiple hypotheses in ambiguous environments but integrates evidence over time to resolve such ambiguity.

Most of the existing work assumes an upright camera pose [16, 17, 20, 28], while some methods [16, 17] explicitly only consider panorama images. In contrast, our method is designed to work with low-cost sensors, *e.g.*, those readily available in all modern phones. Our framework takes only a single perspective image per time-step but operates at a high speed to allow for the frequent integration of new data. To cope with poses with non-zero roll-pitch angle, we utilize the data of an inertial measurement unit and propose a novel data augmentation method to overcome the limitation of previous methods [16, 17, 20, 28].

In this paper we propose the following contributions. **i**) We base our model on a novel 1D ray representation that reflects the 2D floorplan representation. **ii**) We extract scene geometry from single and multi-view cues. A novel selection network fuses them in dependence of the current relative poses to take advantage of either methodology. **iii**) A data augmentation technique using virtual roll-pitch overcomes the limitations of current state-of-the-art methods and allows to cope with non-zero roll-pitch angles in practical use cases. **iv**) To eliminate ambiguity and boost localization, the predictions are filtered over time by a novel and efficient histogram filter formulated as grouped convolution from ego-motion. **v**) Our full system outperforms the state-of-the-art methods in both accuracy and efficiency on existing benchmarks and a real world experiment further illustrates its potential for practical applications. **vi**) We collect a large indoor dataset, composed of floorplans and both short and long sequential observations in 119 Gibson [41] indoor environments. The dataset will be released publicly.

## 2. Related Work

**Visual Localization** is one of the oldest problems in computer vision and is addressed by using various methodologies. Image retrieval based methods [40][1][2] find the most similar image in a database and estimate the query image pose using the pose of the retrieved one. Methods based on a pre-built 3D SfM model of the environment [24, 34, 37–39] establish 2D-3D correspondences between a query image and the 3D structure by matching local descriptors and compute the image pose using minimal solvers and RANSAC.

Recent data driven models deviate from these classical pipelines. Scene coordinate regression [6][42][44] learns to regress the 3D coordinates of the pixels in the query image. Pose regression methods [21][47][50] use a neural network to directly regress a 6D camera pose from the input image. These methods rely on a pre-built 3D model that requires large storage and are scene-specific, which renders them unable to handle unvisited environments.

Instead of using a 3D model to recover the full 6D camera pose, some works tackle localization with an overhead image, such as a map [33, 36], a satellite patch[51, 56] or a floorplan[16, 17, 28] to estimate the SE2 camera pose or R2 camera location. These methods can localize in unvisited scenes as long as some form of map is provided.

**Floorplan localization** is often associated with Lidar localization [3, 4, 23, 27, 48]. However, the use of Lidar inhibits the usability on common mobile devices. Similar geometric cues can be obtained from other sources, such as point cloud reconstruction from a depth camera [18] or Visual Odometry (VO) [8]. [5] extract room edges and compare them against the floorplan layout. To reconstruct 3D geometry, these works usually assume the knowledge of room or camera height [5, 8]. Recently, learning-based methods use only RGB images to localize in a floorplan. LaLaLoc [17] estimates the position of a panorama image in a given floorplan. Assuming known camera and ceiling height, panoramic depth images are rendered at sampled positions within the floorplan. Localization is achieved by comparing map and image features that are embedded into the same feature space during training. LaLaLoc++ [16] eliminates the assumption of known camera and ceiling height by directly embedding the entire floorplan into the feature space. Laser [28] represents the floorplan as a set of points and gathers features, embedded by Pointnet [30], of the visible points for each pose in the floorplan. Images are embedded into a circular feature lying in the same space as the pose features. Similar to LaLaLoc and LASER, our
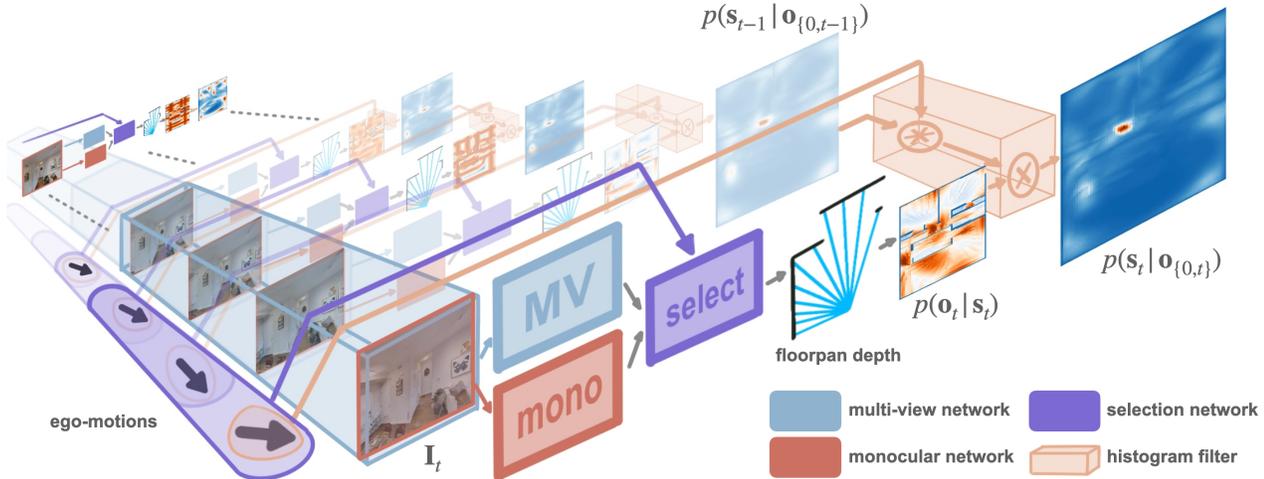
Figure 2. **Pipeline overview.** Our pipeline adopts a monocular (Sec. 3.2) and a multi-view network (Sec. 3.3) to predict floorplan depth. A selection network (Sec. 3.4) consolidates both predictions based on the relative poses. The resulting floorplan depth is used in our observation model and integrated over time by our novel SE(2) histogram filter (Sec. 3.5) to perform sequential floorplan localization.

framework actively compares rendered pose features and query image features to localize.

PF-net [20] tackles visual floorplan localization within a differentiable particle filtering framework. Its observation model is a learned similarity between the image and the corresponding front-facing map patch. The entire system is end-to-end trainable. However, their observation model does not appear as strong as those in [16, 17, 28].

[16, 17, 20, 28] all assume that the images are captured with an upright camera pose. This is a strong requirement for devices such as head-mounted or hand-held devices and appears impractical for some VR/AR use cases. Particularly, LaLaLoc [17] and LaLaLoc++ [16] only work with panorama images, restricting their deployment on most mobile devices. In contrast, we propose a data augmentation scheme to cope with non-upright camera poses, improving the practicability of the method. Furthermore, our method utilizes 1D-range images as internal representation, instead of unorganized point cloud data, 2D-depth or RGB images.

LASER [28] and SeDAR [27] use semantic information such as windows and doors as additional source of information. Because such data is not always present in any floorplan, we consider only occupancy information in this work.

**Sequential localization** , *i.e.* integrating predictions over time can increase the robustness against the observation model, eliminate scene ambiguities and boost the performance of localization [35, 36]. A common framework for fusing sequential observations is the Bayesian filter [5, 8, 10, 18, 20, 27], which maintains the posterior distribution of the current pose in an online fashion. Implementations differ in the representation of the posterior, which can be Gaussian belief (Kalman Filter [49]), a histogram (Histogram

Filter [19]) or weighted particles (Particle Filter [45]). As mentioned, PF-net [20] introduces the particle filter specifically for floormap localization. Here we argue that a histogram filter allows for more scaleable and effective filtering. [19] consider the measurement update as elementwise multiplication, and transition as convolution. However, the presented 1D and 2D cases are not practical for our localization tasks that require at least SE2 pose estimation. To this end, we propose to consider the SE2 motion update as grouped convolution with transition filters derived from known ego-motion that can be implemented efficiently.

**Depth Estimation** provides strong geometric information for localization. Recent advances in deep learning have enabled dense depth prediction from a single image [12–14, 31, 32, 43]. However, monocular depth estimation can suffer from scale ambiguity. In contrast, given sufficient baseline between views, Multi-view stereo (MVS) [9] does not suffer from this problem. Current data-driven MVS methods [7, 25, 26, 29, 52, 53] use neural networks to extract features and learn to filter a cost volume. Instead of estimating pixel-wise depth we predict the floorplan depth of each column of the most recent gravity-aligned image, which can be compared directly against the floorplan. Moreover, we benefit from the advantages of either technique by learning to fuse their predictions.

## 3. Method

### 3.1. Problem Definition and Overview

We solve the problem of localizing RGB images with respect to a floorplan. Given a temporal sequence of $k + 1$ RGB images $\mathcal{I} = \{\mathbf{I}_\tau | \tau \in \{t - k, \cdots, t\}\}$ with known rel-

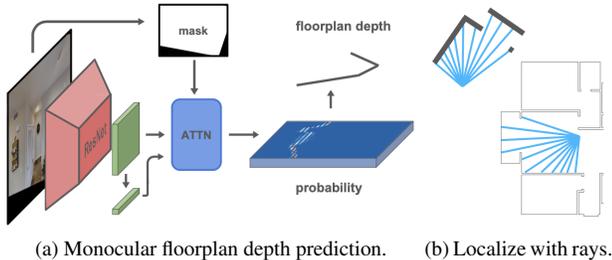(a) Monocular floorplan depth prediction.  (b) Localize with rays.

Figure 3. **Predicting and localizing with a single image.** (a) A gravity aligned image is fed into the ResNet [15] and Attention [46] based feature network. Invisible pixels are masked out in the attention. The network outputs a probability distribution over depth hypotheses and its expectation is used as predicted floorplan depth. (b) Equiangular rays are interpolated from the predicted floorplan depth. We localize by finding the pose in the floorplan that has the most similar rays as the prediction.
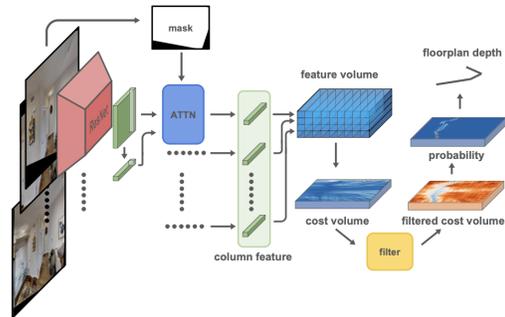


Figure 4. **Floorplan depth prediction from multiple views.** Column features of the images are extracted and gathered in the reference frame. Their cross-view feature variance is used as cost. A U-Net-like network learns the cost filtering to form a probability distribution, and the floorplan depth is defined by its expectation.

ative poses, camera intrinsics, and gravity directions, we aim to find the current $SE(2)$ camera pose $\mathbf{s}_t$ within a given 2D floorplan, where $\mathbf{s}_t = [s_{x,t}, s_{y,t}, s_{\phi,t}]$ represents the camera $x, y$ coordinate in the floorplan and its orientation. We assume the floorplan to encode necessary geometric occupancy information such as doors and walls but no semantic classes. An example is illustrated in Fig. 3 b.

We first estimate the floorplan depth (*i.e.*, the depth to the floorplan occupancy) from the current (Sec. 3.2) and a few recent frames with known relative poses (Sec. 3.3). An MLP fuses the two estimations based on the relative poses and their respective mean depth prediction (Sec. 3.4). We interpolate equiangular rays from the floorplan depth before using them to localize within the floorplan. A histogram filter efficiently fuses the current with integrated past belief, through grouped convolution (Sec. 3.5) to deliver the final localization. The pipeline is illustrated in Fig. 2.

### 3.2. Single Image Localization

We first align the image with the gravity direction and use a ResNet[15] and Attention [46] based network to learn a probability distribution of the floorplan depth over a range of depth hypotheses. Pixels that become unobservable by the gravity alignment are masked out in the attention. The expectation is used as the floorplan depth prediction as illustrated in Fig. 3 a. Finally, we construct an equiangular ray scan from the predicted floorplan depth to localize in the floorplan, compare Fig. 3 b. The more compact representation renders the descriptor independent of the acquisition device and allows for the offline construction of the map pose features, i.e., via a circular equiangular ray scan.

### 3.3. Multiview Stereo Estimation

Inspired by multiview stereo, we adopt a variant of the MVS network [52, 55] to estimate the floorplan depth from mul-

tiple frames with known relative poses. We first extract features of the image columns using a ResNet [15] and Attention [46] based network, and a gravity alignment mask is used in the attention. With multiple depth hypotheses, the column features from different views are gathered via plane sweeping into the reference frame. This procedure is commonly used in dense multiview depth prediction [52] with the exception that we reduce our depth prediction and features vertically instead of predicting depth and extracting features for every pixel. Details can be found in the supplementary material.

The cross-view feature variance forms a cost distribution over the depth hypothesis. We incorporate the observability of the features at different depth hypotheses to compute meaningful variance. Unlike traditional multiview stereo methods [7, 25, 26, 29, 52, 53] that construct 3D (without the channel dimension) cost volumes, we yield 2D cost distribution. As a consequence, the learned cost filter is 2D convolution instead of 3D. A soft-argmin computes the final floorplan depth from the filtered cost distribution as

$$d = \mathbf{d}_{\text{hyp}}^{\top}\text{softmax}(-\mathbf{c}), \qquad (1)$$

where $\mathbf{d}_{\text{hyp}} \in \mathbb{R}^D$ is the vector containing the $D$ depth hypotheses, $\mathbf{c} \in \mathbb{R}^D$ is the cost at each hypothesis, and $\text{softmax}(-\mathbf{c})$ is the probability of each hypothesis.

### 3.4. Learned Complementary Selection

While monocular depth estimation is independent of camera motion but prone to scale ambiguity, Multiview stereo approaches [52, 55] deliver correct scale, but rely on sufficient baselines and camera overlap. Based on these observations, we adopt another MLP that softly selects from the two predictions. The network takes the relative poses of the frames and the estimated multiview and monocular mean floorplan depth as inputs and outputs the correspond-

ing weight for the two estimates. The probability distributions are then fused as the weighted average, *i.e.*,

$$\mathbf{P}_{\text{fuse}} = w\mathbf{P}_{\text{mono}} + (1-w)\mathbf{P}_{\text{mv}}, \quad (2)$$

where $0 \le w \le 1$ is the output by the MLP, $\mathbf{P}_{\text{mono}}$ and $\mathbf{P}_{\text{mv}}$ denote the probability distributions from a single and multi view, respectively. The expectation of the fused probability distribution $\mathbf{P}_{\text{fuse}}$ then provides the final depth prediction.

### 3.5. Sequential Localization

We use a histogram filter to keep track of the posterior over the entire floorplan. We use the predicted floorplan depth as our observation and the following observation model

$$p(\mathbf{s}_t | \mathbf{o}_t) = e^{-\|\hat{\mathbf{r}} - \mathbf{r}_{\mathbf{s}_t}\|_1}, \quad (3)$$

where $\mathbf{r}_{\mathbf{s}_t}$ is the floorplan ray at pose $\mathbf{s}_t$ and $\hat{\mathbf{r}}$ is the interpolated ray from the floorplan depth prediction. We use the relative pose between frames, *i.e.*, ego-motion as the transition model

$$\mathbf{s}_{t+1} = \mathbf{s}_t \oplus \mathbf{t}_t + \omega_t, \quad (4)$$

where $\mathbf{t}_t = [t_{x,t}, t_{y,t}, t_{\phi,t}]$ and $\omega_t = [\omega_{x,t}, \omega_{y,t}, \omega_{\phi,t}]$ are the ego-motion and transition noise at time $t$, respectively, the operator $\oplus$ applies an ego-motion on a state. Further assuming the transition noise $\omega_t$ obeys a Gaussian distribution, the transition probability is expressed as

$$p(\mathbf{s}_{t+1}|\mathbf{s}_t, \mathbf{t}_t) = e^{-\frac{1}{2}(\mathbf{s}_{t+1} - \mathbf{s}_t \oplus \mathbf{t}_t)^\top \Sigma^{-1}(\mathbf{s}_{t+1} - \mathbf{s}_t \oplus \mathbf{t}_t)}, \quad (5)$$
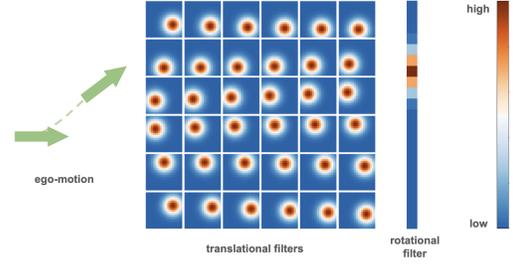
where we model $\Sigma = \text{diag}(\sigma_x^2, \sigma_y^2, \sigma_\phi^2)$ as covariance of the Gaussian distribution. Applying Bayes rule yields

$$p(\mathbf{s}_{t+1}|\mathbf{o}_t, \mathbf{t}_t) = \frac{1}{Z} \sum_{\mathbf{s}_t} p(\mathbf{s}_{t+1}|\mathbf{s}_t, \mathbf{t}_t) p(\mathbf{s}_t|\mathbf{o}_t), \quad (6)$$
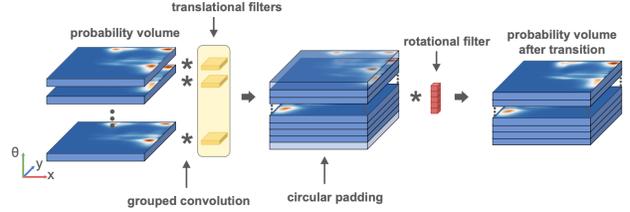
where $Z$ is a normalization factor.

In the following we drop the subscripts indicating the time-step for simplicity. Our histogram filter represents the posterior as a 3D probability volume containing the probability of being at pose $[s_x, s_y, s_\phi]$. The transition is implemented as transition filters [19]. Unlike previous work [19] operating on euclidean $\mathbb{R}^2$ state space, we work on SE2 and transit through ego-motions, so the translation in the world frame depends on the current orientation. Therefore, we decouple the translation and rotation and apply different 2D translation filters for different orientations, before applying the rotation filter to the entire volume along the orientation axis. The 2D translation step can be implemented efficiently as a grouped convolution [22], where each orientation is a group as illustrated in Fig. 5. The translational filter $\mathbf{T}_\phi$ for orientation $\phi$ can be computed through

$$\mathbf{T}_\phi(x, y) = e^{-\frac{1}{2}\delta\mathbf{t}^\top \text{diag}(\sigma_x^2, \sigma_y^2)^{-1}\delta\mathbf{t}}, \quad (7)$$



(a) Transitional filters.



(b) Transition.

Figure 5. **Transition as grouped convolution.** (a) Illustration of the translational filters (from left to right, top to bottom the filters for 0, 10 to 350° ) and the rotational filters derived from a sample ego-motion. (b) The probability volume is divided into $O$ groups, where $O$ is the number of orientations. Each group is convolved with its respective translational filter and stacked back together. After circular padding along the orientation axis, the volume is convolved with the rotational filter to finish the transition step.

where

$$\delta\mathbf{t} = \mathbf{R}_{s_\phi}^{-1}[s_x, s_y]^\top - [t_x, t_y]^\top \quad (8)$$

with $\mathbf{R}_{s_\phi} \in \mathbb{R}^{2\times2}$ being the rotation matrix with angle $s_\phi$. The rotational filter $\mathbf{r}$ is

$$\mathbf{r}(\phi) = e^{-\frac{1}{2}(\phi - t_\phi)^2/\sigma_\phi^2}. \quad (9)$$

The pose posterior corresponds to the filtered probability volume and we can obtain the (best) pose prediction and its uncertainty by a lookup.

## 4. Training

### 4.1. Dataset

We collect a customized dataset with perspective images of 108° horizontal field of view in iGibson [41], an indoor simulation environment, and manually label the floorplans (see Fig. 6) from the provided mesh. The dataset consists of 118 distinct indoor environments and is partitioned into training (100), validation (9), and test (9) sets. We collected three datasets according to the type of motions designed to be typical trajectories for a human holding a phone. One including in-place turning, which we refer to as Gibson(g) for general motions, containing 49558 pieces of 4 sequential views, one without (Gibson(f) for forward motions), containing 24779 pieces of 4 sequential views,

and one containing 118 pieces of 280 to 5152 steps long trajectories(Gibson(t) for trajectories). We also evaluate the proposed single frame localization on Structured3D [54], a photorealistic dataset containing 3296 fully furnished indoor environments with in total 78453 perspective images with 80° horizontal field of view. For Structured3D, we follow the official split.

## 4.2. Virtual Roll Pitch Augmentation

To cope with non-upright camera poses we propose an augmentation technique during the training through virtual roll pitch angle simulation. Perspective images with the same principle point and different viewing angles relate to each other through a simple homography as shown in Fig. 8a. With known camera intrinsic matrix $\mathbf{K}$, camera roll and pitch angle $\psi$, $\theta$, the homography from the original image to the gravity-aligned image is

$$\hat{\mathbf{p}} = \mathbf{K}\mathbf{R}\mathbf{K}^{-1}\mathbf{p}, \qquad (10)$$

where $\mathbf{p}, \hat{\mathbf{p}}$ are the homogeneous image coordinates of the original pixel and the corresponding pixel in the gravity-aligned image, $\mathbf{R}$ is the rotation matrix to the gravity-aligned pose. To simulate the virtual roll pitch angle, we use this to calculate which pixel is observable at angle $\psi$, $\theta$ and mask out the unobservable ones. This is equivalent to the gravity-alignment of the image taken at angle $\psi$, $\theta$.

## 4.3. Training Scheme

Details on the training procedure can be found in the supplementary material. For all training, we optimize the L1 loss to the ground truth floorplan, except for the monocular network, for which we added a shape loss computed as the cosine similarity, *i.e.*,

$$\mathcal{L} = ||\mathbf{d}, \mathbf{d}^*||_1 + \lambda \frac{\mathbf{d}^\top \mathbf{d}^*}{\max\{||\mathbf{d}||_2 ||\mathbf{d}^*||_2, \epsilon\}}, \qquad (11)$$

where $\mathbf{d}, \mathbf{d}^*$ are the predicted and the ground truth depth and $\epsilon$ a small constant to prevent from division by zero.

## 5. Results

We compare our method with the state-of-the-art floorplan localization methods PF-net [20] and LASER [28], both without semantic labels. We sample pose position and orientation at a resolution of 0.1m×0.1m and 10°.

### 5.1. Observation Model

Fig. 6 provides a qualitative comparison between the methods. While all predictions possess multiple modes, our probability estimate appears more accurate, due the accurate floorplan depth estimation and the invariance of the ray representation. In the following, we thoroughly investigate the performance of the proposed observation model.

**Single Frame**. We evaluate single frame localization accuracy on Gibson(f) and Structured3D. As shown in Tab. 1, the proposed monocular network, Ours$_s$ significantly outperforms both baselines on Gibson(f) seeing almost 200% improvement across all metrics. Also on Structured3D our method surpasses the state of the art by a large margin. When taking the orientation into account, the recall does not drop much (35.1% at 1m 30 deg compared to 36.6% at 1m on Gibson(f) and 21.3% to 22.4% on Structured3D). This underlines the accurate orientation estimation of our method. We notice here the performance of LASER on Structured3D does not align with that reported in [28], we suspect this is due to slightly difference in the dataset (*perspective Structured3D* compared to their perspective images cropped from *panoramic Structured3D*) and the random roll pitch angles this dataset contains.

**Multiview**. Because the existing indoor datasets either do not provide sequential images or a floorplan, we evaluate the proposed multiview module, Ours$_m$ only on the collected Gibson dataset. Tab. 1 verifies that the multiview module can clearly outperform the two baselines on Gibson(f), and notably increases the recall by more than 20% at all thresholds compared to our monocular module. This shows the effectiveness of using multiview geometry cues in the observation model for floorplan depth estimation.

However, multiview estimation can suffer from small baselines or insufficient overlap present in the general mo-

| R@ | Gibson(f) | | | | Structured3D | | | |
|---|---|---|---|---|---|---|---|---|
| | 0.1m | 0.5m | 1m | 1m30° | 0.1m | 0.5m | 1m | 1m30° |
| PF-net | 0 | 2.0 | 6.9 | 1.2 | 0.2 | 1.3 | 3.2 | 0.9 |
| LASER | 0.4 | 6.7 | 13.0 | 10.4 | 0.7 | 6.4 | 10.4 | 8.7 |
| Ours$_s$ | **4.7** | **28.6** | **36.6** | **35.1** | **1.5** | **14.6** | **22.4** | **21.3** |
| Ours$_m$ | **13.2** | **40.9** | **45.2** | **43.7** | - | - | - | - |

Table 1. **Comparison between our observation model and the baselines.**

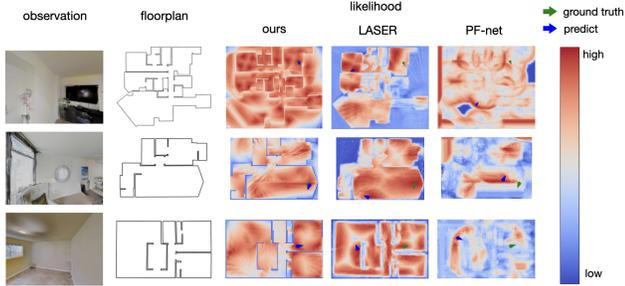| R@ | Gibson(g) | | | |
|---|---|---|---|---|
| | 0.1m | 0.5m | 1m | 1m30° |
| PF-net | 1.0 | 1.9 | 5.6 | 1.9 |
| LASER | 0.7 | 7.0 | 11.8 | 9.5 |
| Ours$_s$ | 4.3 | 26.7 | 33.7 | 32.3 |
| Ours$_m$ | 9.3 | 27.0 | 31.0 | 29.2 |
| Ours$_t$ | 10.5 | 34.3 | 39.6 | 38.0 |
| Ours$_f$ | **12.2** | **39.4** | **44.5** | **43.2** |

Table 2. **Complement single and multiview.**

Figure 6. **Single observation likelihood.** Utilizing the front-facing map patch, PF-net does not account for occlusion or the camera's field of view. Using a set of point features, LASER is not invariant to rotation and translation. Contrary, our 1D ray-scan representation possess such invariance and inherently considers occlusions.
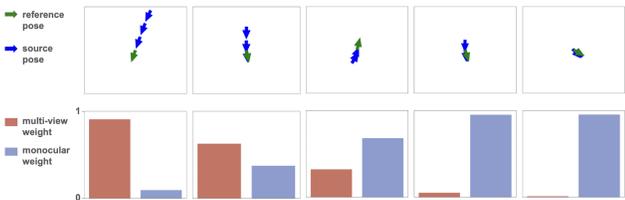


Figure 7. **Selection weights VS relative poses.** The section network computes a weighted combination of monocular and multi-view observation estimate. The more degenerate the relative poses become the more the monocular estimate is preferred.

tion dataset, Gibson(g), which includes nearly in-place rotation. Here, as shown in Tab. 2, the recall of the multiview module falls below that of the monocular module for both larger thresholds, 1m and 1m 30°.

**The selection network** is evaluated on the general motion dataset Gibson(g) in Tab. 2. The selection network, $Ours_f$, delivers a 30-50% improvement across all precisions compared to both individual networks $Ours_s$ and $Ours_m$. As a baseline selection we also evaluate selection by thresholding the relative motions between sequential frames named $Ours_t$. While this baseline achieves a 18% improvement over the individual networks, which further proves the idea of complementing monocular with multiview estimation, the selection network learns a more sophisticated selection rule and achieves at least an additional extra 12% improvement. Examples of the selection decisions are illustrated in Fig. 7.

**Virtual Roll-Pitch.** Fig. 8b compares the recall of the monocular module trained with and without virtual roll-pitch augmentation on Gibson(f) at 1m×1m×30° resolution. The network trained without augmentation shows decreasing recall when the roll pitch disturbances are imposed (especially for large pitch angles). In contrast the proposed virtual roll-pitch augmentation increases the robustness of



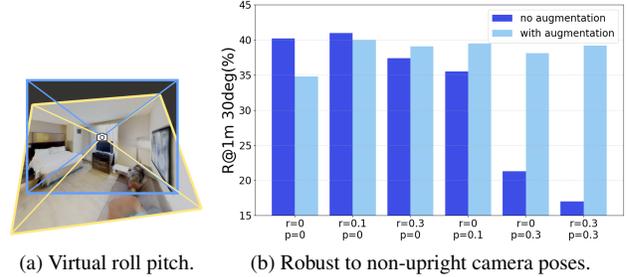(a) Virtual roll pitch.    (b) Robust to non-upright camera poses.

Figure 8. **Virtual roll pitch augmentation.** (a) After gravity alignment we mask out unobservable pixels (in black). During training we augment the data accordingly. (b) If trained without augmentation, the recall of the network decreases as the roll and pitch angle increases. Training with augmentation significantly increases robustness against non-upright camera poses.
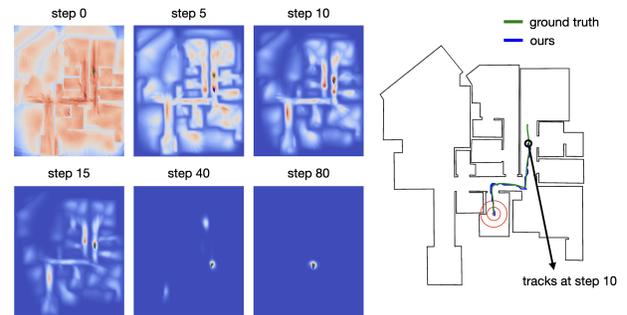


Figure 9. **Posterior evolution and trajectory.** Our strong observation model already provides an accurate estimation at the initial step. Due to the ambiguous nature of the floorplan (the hallway), the posterior estimates shows multi-modality. After 10 steps, our system tracks firmly at a frequency of 27Hz on this 18.4m×15.5m floorplan using a laptop NVidia RTX 3070Ti GPU.

the recall against non-upright poses.

## 5.2. Sequential localization

Our full sequential localization pipeline is evaluated on the Gibson(t) dataset, containing long simulated trajectories. A qualitative study in Fig. 9 shows that the proposed histogram filter can effectively maintain a global posterior of the camera pose. At the start the distribution has multiple modes, as the camera movement provides more and more evidence, the distribution converges to a single sharp peak.

| | LASER | $Ours_s$ | $Ours_f$ |
|---|---|---|---|
| Success rate@1m (%) | 59.5 | 89.2 | **94.6** |
| RMSE(succeeded) (m) | 0.39 | 0.18 | **0.12** |
| RMSE(all) (m) | 1.96 | 0.88 | **0.51** |

Table 3. **Comparison of observation models integrated into our histogram filter.** RMSEs are computed from the last 10 frames
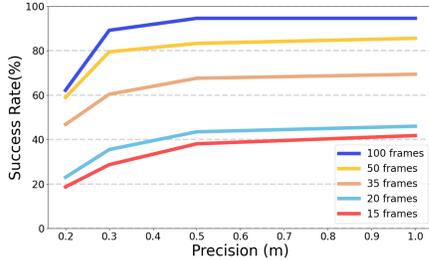
Figure 10. **Localization success rate vs. precision threshold for different filter history sizes**. The more frames are used within the filter, the higher the localization success rate.

**Success Rate**. As a metric we consider sequential localization at Xm as successful if the prediction stays within a radius of Xm over the last 10 frames. We integrate the baseline observation models into our histogram filter and compare them against our pipeline in Tab. 3. Our full system achieves a success rate of 94.6% at 1m using a history of 100 frames, surpassing the two baselines by more than 58%, and our monocular observation by 10%. We also compare the RMSE (over the last 10 frames) of our trajectory tracking in both *succeeded* and *all* runs. Here, our full pipeline delivers 70% lower error (0.12 m and 0.51 m, respectively) than the baselines. We compare success rates for various number of frames in Fig. 10. In general, using more frames increases the success rates.

**Timing.** Tab. 4 compares the runtime of different combinations of observation and filtering models. Despite the slightly slower feature extraction of our proposed observation model, the rapid matching helps it to achieve the highest iteration rates. The particle filter (PF) suffers from expensive resampling and feature rendering and demands instanciating a large number of samples for global localization in a large area. Analogously, our histogram filter (HF) utilizes presampled "particles", constructed offline, and can avoid constant rerendering at runtime. As a result, our histogram filter achieves 45% faster iteration than the particle filter.

## 6. Real-world Experiment

Since no real-world indoor dataset with both sequential observations and floorplan exist that allows training and test-

|  | Feature Extraction(s) | Matching(s) | Iteration | |
|---|---|---|---|---|
|  |  |  | HF | PF |
| PF-net(obs) | 0.042 | 2.375 | - | - |
| LASER(obs) | 0.008 | 0.224 | 0.241 | 0.287 |
| Ours$_f$ | 0.033 | 0.003 | 0.037 | 0.067 |

Table 4. **Timing.** Because PF-net is too slow we do not test its performance in filters.
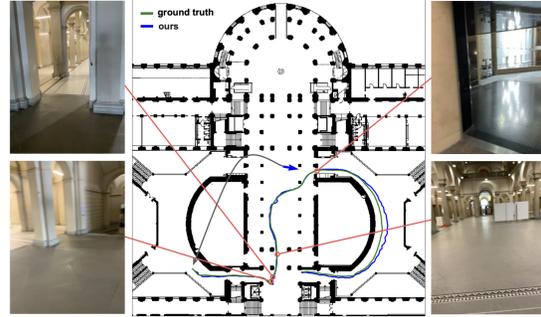


Figure 11. **Sequential localization in HGE.** The localization area is 75m×81m with challenging observations including motion blur, non-lambertian surfaces, ambiguities and occlusions. Our trajectory tracks the ground truth closely from the second step. It deviates slightly later due to the ambiguous floorplan labeling, however, recovers shortly thanks to the filters and converges to a sharp posterior estimation in the end.

ing, we show the potential of our pipeline in real-world scenario by customizing LaMAR [35]. LaMAR is a real world dataset containing three scenes. We select the trajectories in HGE indoor scene containing trajectories within a single floor, and split it into training and testing set. We use our single frame observation model with the proposed histogram filter to localize. The entire floor has an area of 80m×120m, and the data includes challenging observations as shown in Fig. 11. We use the data within 75m×81m and localize within the floorplan of the same size. Our system localizes and tracks the camera pose from the second step and closely follows it afterwards. Despite the large scene scale, our histogram filter is still efficient enough to localize at 3 hz.

## 7. Limitations and Conclusion

Through the process, we realized a lack of indoor datasets with sequential observations and floorplan. Although we tried to mitigate this by collecting a dataset in a simulated environment, more real-world datasets are highly desirable to close the domain gap. While our proposed system effectively uses geometric cues, ambiguities could be further reduced by utilizing semantic information from both the image and the floorplan. In this work, we present a data-driven and probabilistic model for localization within a floorplan. The system is more practical than previous methods, demanding only consumer hardware, perspective RGB images and non-upright camera poses, while operating at very high frame-rates. Our system allows for both accurate single-frame and sequential localization in unvisited environments. It outperforms the state-of-the-art in both tasks across different datasets and various metrics by a significant margin. Finally, we illustrate its real world potential on a challenging large scale indoor dataset. Our work could be interesting in many indoor AR/VR applications and boost robot auton-

omy in indoor environments.

# References

[1] Relja Arandjelovic, Petr Gronat, Akihiko Torii, Tomas Pajdla, and Josef Sivic. Netvlad: Cnn architecture for weakly supervised place recognitio. In *CVPR*, pages 5297–5307, 2016. 1, 2

[2] Vassileios Balntas, Shuda Li, and Victor Prisacariu. Relocnet: Continuous metric learning relocalisation using neural nets. In *ECCV*, pages 751–767, 2018. 1, 2

[3] Federico Boniardi, Tim Caselitz, Rainer Kümmerle, and Wolfram Burgard. Robust lidar-based localization in architectural floor plans. In *IROS*, pages 3318–3324, 2017. 2

[4] Federico Boniardi, Tim Caselitz, Rainer Kümmerle, and Wolfram Burgard. A pose graph-based localization system for long-term navigation in cad floor plans. pages 84–97, 2019. 2

[5] Federico Boniardi, Abhinav Valada, Rohit Mohan, Tim Caselitz, and Wolfram Burgard. Robot localization in floor plans using a room layout edge extraction network. In *IROS*, pages 5291–5297, 2019. 2, 3

[6] Eric Brachmann, Alexander Krull, Sebastian Nowozin, Jamie Shotton, Frank Michel, Stefan Gumhold, and Carsten Rother. Dsac-differentiable ransac for camera localization. In *CVPR*, pages 6684–6692, 2017. 2

[7] Shuo Cheng, Zexiang Xu, Shilin Zhu, Zhuwen Li, Li Erran Li, Ravi Ramamoorthi, and Hao Su. Deep stereo using adaptive thin volume representation with uncertainty awareness. In *CVPR*, pages 2524–2534, 2020. 3, 4

[8] Hang Chu, Dong Ki Kim, and Tsuhan Chen. You are here: Mimicking the human thinking process in reading floorplans. In *ICCV*, pages 2210–2218, 2015. 2, 3

[9] Robert T Collins. A space-sweep approach to true multi-image matching. In *CVPR*, pages 358–363, 1996. 3

[10] Frank Dellaert, Dieter Fox, Wolfram Burgard, and Sebastian Thrun. Monte carlo localization for mobile robots. In *ICRA*, pages 1322–1328, 1999. 3

[11] Jeffrey Delmerico, Stefano Mintchev, Alessandro Giusti, Boris Gromov, Kamilo Melo, Tomislav Horvat, Cesar Cadena, Marco Hutter, Auke Ijspeert, Dario Floreano, et al. The current state and future outlook of rescue robotics. *Journal of Field Robotics*, 36(7):1171–1191, 2019. 2

[12] David Eigen and Rob Fergus. Predicting depth, surface normals and semantic labels with a common multi-scale convolutional architecture. In *NeurIPS*, pages 2650–2658, 2015. 3

[13] David Eigen, Christian Puhrsch, and Rob Fergus. Depth map prediction from a single image using a multi-scale deep network. 2014.

[14] Clément Godard, Oisin Mac Aodha, and Gabriel J Brostow. Unsupervised monocular depth estimation with left-right consistency. In *CVPR*, pages 270–279, 2017. 3

[15] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *CVPR*, pages 770–778, 2016. 4

[16] Henry Howard-Jenkins and Victor Adrian Prisacariu. Lalaloc++: Global floor plan comprehension for layout localisation in unvisited environments. In *ECCV*, pages 693–709, 2022. 1, 2, 3

[17] Henry Howard-Jenkins, Jose-Raul Ruiz-Sarmiento, and Victor Adrian Prisacariu. Lalaloc: Latent layout localisation in dynamic, unvisited environments. In *ICCV*, pages 10107–10116, 2021. 1, 2, 3

[18] Seigo Ito, Felix Endres, Markus Kuderer, Gian Diego Tipaldi, Cyrill Stachniss, and Wolfram Burgard. W-rgb-d: floor-plan-based indoor global localization using a depth camera and wifi. In *ICRA*, pages 417–422, 2014. 2, 3

[19] Rico Jonschkowski and Oliver Brock. End-to-end learnable histogram filters. In *Workshop on Deep Learning for Action and Interaction at NIPS*, 2016. 3, 5

[20] Peter Karkus, David Hsu, and Wee Sun Lee. Particle filter networks with application to visual localization. In *CoRL*, pages 169–178, 2018. 1, 2, 3, 6

[21] Alex Kendall, Matthew Grimes, and Roberto Cipolla. Posenet: A convolutional network for real-time 6-dof camera relocalization. In *ICCV*, pages 2938–2946, 2015. 2

[22] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. In *NeurIPS*, 2012. 5

[23] Zhikai Li, Marcelo H Ang, and Daniela Rus. Online localization with imprecise floor space maps using stochastic gradient descent. In *IROS*, pages 8571–8578. 2

[24] Liu Liu, Hongdong Li, and Yuchao Dai. Efficient global 2d-3d matching for camera localization in a large-scale 3d map. In *ICCV*, pages 2372–2381, 2017. 1, 2

[25] Keyang Luo, Tao Guan, Lili Ju, Haipeng Huang, and Yawei Luo. P-mvsnet: Learning patch-wise matching confidence aggregation for multi-view stereo. In *ICCV*, pages 10452–10461, 2019. 3, 4

[26] Keyang Luo, Tao Guan, Lili Ju, Yuesong Wang, Zhuo Chen, and Yawei Luo. Attention-aware multi-view stereo. In *CVPR*, pages 1590–1599, 2020. 3, 4

[27] Oscar Mendez, Simon Hadfield, Nicolas Pugeault, and Richard Bowden. Sedar: Reading floorplans like a human—using deep learning to enable human-inspired localisation. *IJCV*, 128:1286–1310, 2020. 2, 3

[28] Zhixiang Min, Naji Khosravan, Zachary Bessinger, Manjunath Narayana, Sing Bing Kang, Enrique Dunn, and Ivaylo Boyadzhiev. Laser: Latent space rendering for 2d visual localization. In *CVPR*, pages 11122–11131, 2022. 1, 2, 3, 6

[29] Rui Peng, Rongjie Wang, Zhenyu Wang, Yawen Lai, and Ronggang Wang. Rethinking depth estimation for multi-view stereo: A unified representation. In *CVPR*, pages 8645–8654, 2022. 3, 4

[30] Charles R Qi, Hao Su, Kaichun Mo, and Leonidas J Guibas. PointNet: Deep learning on point sets for 3d classification and segmentation. In *CVPR*, pages 652–660, 2017. 2

[31] René Ranftl, Katrin Lasinger, David Hafner, Konrad Schindler, and Vladlen Koltun. Towards robust monocular depth estimation: Mixing datasets for zero-shot cross-dataset transfer. *IEEE TPAMI*, 44(3):1623–1637, 2020. 3

[32] René Ranftl, Alexey Bochkovskiy, and Vladlen Koltun. Vision transformers for dense prediction. In *ICCV*, pages 12179–12188, 2021. 3

[33] Noe Samano, Mengjie Zhou, and Andrew Calway. You are here: Geolocation by embedding maps and images. In *ECCV*, pages 502–518, 2020. 2

[34] Paul-Edouard Sarlin, Cesar Cadena, Roland Siegwart, and Marcin Dymczyk. From coarse to fine: Robust hierarchical localization at large scale. In *CVPR*, pages 12716–12725, 2019. 1, 2

[35] Paul-Edouard Sarlin, Mihai Dusmanu, Johannes L Schönberger, Pablo Speciale, Lukas Gruber, Viktor Larsson, Ondrej Miksik, and Marc Pollefeys. Lamar: Benchmarking localization and mapping for augmented reality. In *ECCV*, pages 686–704, 2022. 2, 3, 8

[36] Paul-Edouard Sarlin, Daniel DeTone, Tsun-Yi Yang, Armen Avetisyan, Julian Straub, Tomasz Malisiewicz, Samuel Rota Bulò, Richard Newcombe, Peter Kontschieder, and Vasileios Balntas. Orienternet: Visual localization in 2d public maps with neural matching. In *CVPR*, pages 21632–21642, 2023. 2, 3

[37] Torsten Sattler, Bastian Leibe, and Leif Kobbelt. Fast image-based localization using direct 2d-to-3d matching. In *ICCV*, pages 667–674, 2011. 1, 2

[38] Torsten Sattler, Bastian Leibe, and Leif Kobbelt. Improving image-based localization by active correspondence search. In *ECCV*, pages 752–765, 2012.

[39] Torsten Sattler, Bastian Leibe, and Leif Kobbelt. Efficient & effective prioritized matching for large-scale image-based localization. *PAMI*, 39(9):1744–1756, 2016. 1, 2

[40] Grant Schindler, Matthew Brown, and Richard Szeliski. City-scale location recognition. In *CVPR*, pages 1–7, 2007. 1, 2

[41] Bokui Shen, Fei Xia, Chengshu Li, Roberto Martín-Martín, Linxi Fan, Guanzhi Wang, Claudia Pérez-D'Arpino, Shyamal Buch, Sanjana Srivastava, Lyne Tchapmi, et al. igibson 1.0: A simulation environment for interactive tasks in large realistic scenes. In *IROS*. 2, 5

[42] Jamie Shotton, Ben Glocker, Christopher Zach, Shahram Izadi, Antonio Criminisi, and Andrew Fitzgibbon. Scene coordinate regression forests for camera relocalization in rgb-d images. In *CVPR*, pages 2930–2937, 2013. 2

[43] Zachary Teed and Jia Deng. DeepV2D: Video to depth with differentiable structure from motion. In *ICLR*, 2020. 3

[44] Julien Valentin, Matthias Nießner, Jamie Shotton, Andrew Fitzgibbon, Shahram Izadi, and Philip HS Torr. Exploiting uncertainty in regression forests for accurate camera relocalization. In *CVPR*, pages 4400–4408, 2015. 2

[45] Rudolph van der Merwe, Arnaud Doucet, Nando de Freitas, and Eric Wan. The unscented particle filter. In *NeurIPS*, 2000. 2, 3

[46] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. 2017. 4

[47] Florian Walch, Caner Hazirbas, Laura Leal-Taixe, Torsten Sattler, Sebastian Hilsenbeck, and Daniel Cremers. Image-based localization using lstms for structured feature correlation. In *ICCV*, pages 627–637, 2017. 2

[48] Xipeng Wang, Ryan J Marcotte, and Edwin Olson. Glfp: Global localization from a floor plan. In *IROS*, pages 1627–1632, 2019. 2

[49] Greg Welch and Gary Bishop. An introduction to the kalman filter. Technical Report 95-041, University of North Carolina at Chapel Hill, 1995. 2, 3

[50] Jian Wu, Liwei Ma, and Xiaolin Hu. Delving deeper into convolutional neural networks for camera relocalization. In *ICRA*, pages 5644–5651, 2017. 2

[51] Zimin Xia, Olaf Booij, Marco Manfredi, and Julian FP Kooij. Visual cross-view metric localization with dense uncertainty estimates. In *ECCV*, pages 90–106, 2022. 2

[52] Yao Yao, Zixin Luo, Shiwei Li, Tian Fang, and Long Quan. Mvsnet: Depth inference for unstructured multi-view stereo. In *ECCV*, pages 767–783, 2018. 3, 4

[53] Yao Yao, Zixin Luo, Shiwei Li, Tianwei Shen, Tian Fang, and Long Quan. Recurrent mvsnet for high-resolution multi-view stereo depth inference. In *CVPR*, pages 5525–5534, 2019. 3, 4

[54] Jia Zheng, Junfei Zhang, Jing Li, Rui Tang, Shenghua Gao, and Zihan Zhou. Structured3d: A large photo-realistic dataset for structured 3d modeling. In *ECCV*, pages 519–535, 2020. 6

[55] Huizhong Zhou, Benjamin Ummenhofer, and Thomas Brox. Deeptam: Deep tracking and mapping. In *ECCV*, pages 822–838, 2018. 4

[56] Sijie Zhu, Taojiannan Yang, and Chen Chen. Vigor: Cross-view image geo-localization beyond one-to-one retrieval. In *CVPR*, pages 3640–3649, 2021. 2