

LASER: LATent Space REndering for 2D Visual Localization

Zhixiang Min¹ Naji Khosravan² Zachary Bessinger² Manjunath Narayana²
 Sing Bing Kang² Enrique Dunn¹ Ivaylo Boyadzhiev²
¹Stevens Institute of Technology ²Zillow Group

Abstract

We present *LASER*, an image-based Monte Carlo Localization (MCL) framework for 2D floor maps. *LASER* introduces the concept of latent space rendering, where 2D pose hypotheses on the floor map are directly rendered into a geometrically-structured latent space by aggregating viewing ray features. Through a tightly coupled rendering codebook scheme, the viewing ray features are dynamically determined at rendering-time based on their geometries (i.e. length, incident-angle), endowing our representation with view-dependent fine-grain variability. Our codebook scheme effectively disentangles feature encoding from rendering, allowing the latent space rendering to run at speeds above 10KHz. Moreover, through metric learning, our geometrically-structured latent space is common to both pose hypotheses and query images with arbitrary field of views. As a result, *LASER* achieves state-of-the-art performance on large-scale indoor localization datasets (i.e. ZInD [5] and Structured3D [38]) for both panorama and perspective image queries, while significantly outperforming existing learning-based methods in speed.

1. Introduction

Camera localization aims to estimate the spatial relationship between a given input image w.r.t. an environmental representation. Diverse application-driven variants have been addressed in the computer vision, robotics, and AR/VR literature. Particular problem instances are defined in terms of the scope of the pose geometric model (e.g. SE(2) vs SE(3)), the type of input query imagery (e.g. RGB, depth), as well as the type of the environmental geometric reference (e.g. geometric maps, registered image collections). Instances where the queries and the geometric reference share the same domain define camera localization as a direct geometric registration problem (e.g. ICP [21], SfM-based geometric verification [26]). Conversely, whenever query observations and the geometric reference are from different domains, it requires the design of integrative cross-modality data representations able to distinguish and associate input observations and reference data.

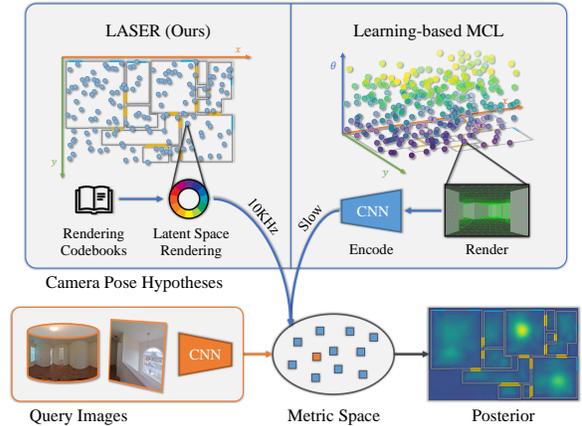


Figure 1. **LASER diagram.** Learning-based MCL frameworks encode camera pose hypotheses and query images into a common metric space to measure their similarities. Compared to existing works, *LASER* directly renders latent features and has a reduced sampling dimension.

This work focuses on solving for the camera pose of a query panorama/perspective image w.r.t. a 2D floor map, under the Monte Carlo Localization (MCL) framework [7]. MCL adopts a generative framework, where the solution (i.e. camera pose) space is systematically sampled to render observation hypotheses and states the problem in terms of a maximum likelihood search and/or optimization w.r.t. a query observation. Note that in this work, we solely focus on improving the measurement model, hence we are only interested in localizing individual queries without initialization. Yet, it is straight-forward to integrate our work into a full MCL framework with customizable temporal updates. Conventional MCL methods [7, 34] require depth sensors and have limited robustness to environmental variability (e.g. furniture/object changes) due to their explicit geometric modeling. Extensions leveraging image-based room layout estimation [2, 32] address content variability at the expense of imposing environmental or capture assumptions such as Manhattan world, known ceiling or camera height. Recent supervised learning approaches [12, 14] have learned a common latent space for both synthesized and

query observations. However, their explicit high-fidelity rendering and CNN-based encoding corresponding to individual pose hypothesis is computationally burdensome. Given the time-sensitivity and the accuracy dependency on the number of samples for MCL applications, their computational burden compromises estimation accuracy for online operation. Moreover, the attained latent space representations are not geometrically interpretable while lacking expressiveness and level of detail due to the coarse-level homogenization common to convolutional architectures. We address these challenges within a geometrically-structured metric learning framework, which performs latent-space rendering while prioritizing applicability to unseen environments, computational efficiency, estimation accuracy and robustness.

We circumvent expensive explicit rendering-and-encoding of sampling observations by directly rendering features in a learnable common metric space from a rasterized 2D floor map. Such latent space rendering is enabled by a rendering codebook, which allows map points to have view-dependent dynamic features for representing rendering-time dynamics (i.e. viewing ray geometries such as length and incident-angle). Importantly, we structure said latent space to be geometrically meaningful by encoding visibility-based omni-directional observations into discretized circular (i.e. angular cyclical) representations. This representation, namely *circular feature*, along with the view-dependent feature encoding from the rendering codebook, provides fine-grain structured descriptors for geometry and semantics at a high sampling FPS of 10KHz. Extensive experiments on Structured3D [38] and ZInD datasets [5] show that our proposed framework significantly outperforms state-of-the-art frameworks both in accuracy and speed. The main technical contributions and innovations driving these performance gains are:

(1) Map-aware 2D visual localization framework: While existing MCL frameworks render hypotheses over a local scope (i.e. contents visible to the camera), LASER uses a 2D variant of the PointNet [22] to attain latent encoding from 2D point cloud maps. The PointNet learns global context from the map and provides the latent feature with map-level scope which improves LASER’s recall.

(2) Latent space rendering based on codebook scheme: LASER obviates the redundant rendering-and-encoding of an intermediate representation for individual samples by directly render features in the latent space. Powered by our rendering codebook scheme, the features are dynamically determined in rendering-time to encode fine-grain ray geometries. Such design achieves significantly higher sampling speed and accuracy at the same time.

(3) Geometrically-structured metric learning: LASER structures the metric learning to be geometrically meaningful using a rotationally-covariant 2D omni-directional cir-

cular feature. Its fine-grain structured variability implicitly expresses environmental layouts for high-accuracy localization, and seamlessly supports query images with arbitrary field of views. In addition, this choice implicitly encodes many orientations that reduces the rotation dimension from the MCL sampling space.

2. Related Works

The general 6-DoF relocation methods either explicitly [3, 19, 23–25, 28, 29, 35, 36, 40] or implicitly [1, 8, 15] find appearance correspondences between the query image and scene representations (e.g. images with known camera poses, sparse/dense 3D reconstructions). The camera pose can then be predicted from neural networks [1, 8, 15] or can be recovered using SfM methods [20, 26, 30, 31, 33, 37, 39] from explicit correspondences. The appearance dependency limits their robustness to appearance changes and cannot work with pure geometry maps (e.g. occupancy map).

Monte Carlo Localization (MCL) [4, 7, 17, 34] is the most popular framework for 2D localization on pure geometry maps. With our interest solely on single query, MCL defines a measurement model over geometry observations from depth sensors and compares it with simulated observations sampled from the floor map. The methods limit the input type to geometric measurements, which also limit their robustness to geometry variation/occlusions in the map (e.g. changes in furniture/object). Some extensions of MCL take intensity images as input. Boniardi *et al.* [2] recover room geometry by explicitly extracting room layout edges using CNNs. Wang *et al.* [32] further cooperates semantic information for localization in large indoor spaces. All these methods are subject to strong assumptions such as Manhattan world and known ceiling or camera height, which limit their applications.

Recent learning-based methods [12–14] extend the MCL into a metric learning framework. They utilize learnable CNNs to encode query images and render location into the same metric latent space to estimate their similarity. For rendering a given camera pose, PfNet [14] uses a spatially transformed bird’s-eye map image, while LaLaLoc [12] assumes known camera and ceiling height, and renders the layout depth image. Given the heavy rendering-and-encoding process, these methods are computationally burdensome, which limits their performance for time-sensitive or SWaP-constrained applications.

Neural Radiance Field (NeRF) [16, 18] is a newly emerging area that synthesizes photo-realistic images from scene-specific neural representations learned using back-propagation. In contrast, our scene representation is the rendering codebooks inferred from PointNet, whose estimation process is scene-agnostic. Moreover, our latent space rendering synthesizes view-dependent latent features.

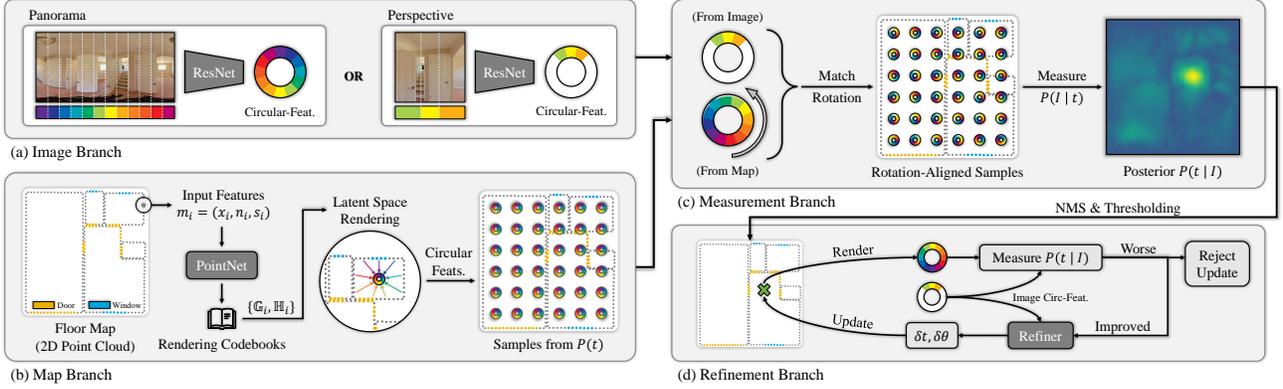


Figure 2. **LASER inference workflow.** Learnable components are shown in dark-gray boxes. (a) Panorama or perspective images are processed into circular features with ResNet50. (b) Rasterized 2D point cloud map is encoded using PointNet into rendering codebooks, from which our uniformly sampled circular features are rendered. (c) The similarities between image and map circular features are measured at their best matched rotation. (d) The final estimation is iteratively refined until its likelihood stops improving.

3. Method

3.1. Problem Formulation

We define camera localization as estimating the 2D pose $\mathbf{p}^* = [\mathbf{t}, \theta] \in SE(2)$, associated with a given query image \mathbf{I} , w.r.t. a reference map \mathbf{M} . The pose parameters $\mathbf{t} \in \mathbb{R}^2$ and $\theta \in [0, 2\pi)$ define, respectively, the camera’s planar displacement vector and yaw axis rotation. The input query \mathbf{I} may be panorama or a perspective image with known FoV. We do not limit the format of map \mathbf{M} , but assume it encodes the occupancy information within a 2D plane.

Monte Carlo Localization. The general Monte Carlo Localization (MCL) framework [7] defines a measurement model $P(\mathbf{I} | \mathbf{p}; \mathbf{M})$, which expresses the likelihood of image \mathbf{I} being observed at camera pose \mathbf{p} on map \mathbf{M} . Henceforth, we obviate the non-random parameter \mathbf{M} from our formulation for simplicity. The posterior distribution of \mathbf{p} after observing \mathbf{I} is the solution of interest. Following Bayes rule, MCL estimates the posterior distribution $P(\mathbf{p} | \mathbf{I})$ as

$$P(\mathbf{p} | \mathbf{I}) = \frac{P(\mathbf{I} | \mathbf{p})P(\mathbf{p})}{P(\mathbf{I})} \quad (1)$$

where $P(\mathbf{I})$ is a normalization constant that can be safely ignored, while $P(\mathbf{p})$ is the prior camera pose distribution, which we assume uniformly distributed within the map area. Finally, the full posterior can be approximated by drawing particles from $P(\mathbf{p})$ whose likelihoods will be estimated using the measurement model as defined in Eq.5.

Localization as Metric Learning. The measurement model $P(\mathbf{I} | \mathbf{p})$ within the MCL framework defines the similarity across the camera pose and image domains. We adopt deep metric learning [11] to learn a unified metric space for comparing cross-domain similarity between the query image and camera pose hypotheses as shown in Fig.2. We

detail how images and camera poses are encoded in to the metric space in §3.3 and §3.2, respectively.

Circular Feature. Contrary to conventional flattened descriptors used in metric learning, we introduce circular features to encode spatial visibility, leading to our geometrically-structured metric learning. We define a circular feature as an ordered set of feature vectors

$$\mathbb{F} = \{\mathbf{f}^\alpha | \alpha = 0 \dots V-1\} \quad (2)$$

where V is the number of feature segments. Each feature segment $\mathbf{f}^\alpha \in \mathbb{R}^D$ encodes a local directional FoV of $\frac{2\pi}{V}$ radians in the range $[\frac{2\pi\alpha}{V}, \frac{2\pi(\alpha+1)}{V})$ on the 2D plane. We denote this ordered set \mathbb{F} as a circular feature since the first and last feature segments correspond to adjacent FoVs. With this design, the omni-directional 2D spatial information is implicitly encoded in the order of feature segments.

Measurement Model. We first define the similarity measurement between two circular features $\mathbb{F}_i = \{\mathbf{f}_i^\alpha | \alpha = 0 \dots V-1\}$ and $\mathbb{F}_j = \{\mathbf{f}_j^\alpha | \alpha = 0 \dots V-1\}$ as

$$\mathcal{S}(\mathbb{F}_i, \mathbb{F}_j) = \frac{\sum_{\alpha=1}^V \cos(\mathbf{f}_i^\alpha, \mathbf{f}_j^\alpha)}{2V} + 0.5 \quad (3)$$

where $\cos(\cdot, \cdot)$ computes the vector cosine similarity, and the function output is normalized to $[0, 1]$. We further define a rotating operator $\mathcal{R}(\mathbb{F}, \theta)$ to rotate the underlying spatial information of a circular feature \mathbb{F} with a given angle θ by

$$\mathcal{R}(\mathbb{F}, \theta) = \{\mathbf{f}^{(\alpha + \frac{V\theta}{2\pi}) \% V} | \alpha = 0 \dots V-1\} \quad (4)$$

where the 1D feature space is linearly interpolated when the indexing yields non integer values. Finally, we define the measurement model as

$$P(\mathbf{I} | \mathbf{p}) = P(\mathbf{I} | \mathbf{t}, \theta) = A \cdot \mathcal{S}(\mathbb{F}_{\mathbf{I}}, \mathcal{R}(\mathbb{F}_{\mathbf{t}}, \theta)) \quad (5)$$

where A is the PDF normalization constant, and $\mathbb{F}_{\mathbf{I}}$ and $\mathbb{F}_{\mathbf{t}}$ are circular features encoded from the query image and rendered at location \mathbf{t} on the map respectively.

Rotation Reduction. As MCL needs a large number of samples to approximate the camera pose posterior in $SE(2)$, we systematically reduce the rotation dimension from the MCL sampling step. For a sample location \mathbf{t} with a canonical orientation, the optimal relative rotation of its circular feature $\mathbb{F}_{\mathbf{t}}$ w.r.t. an image circular feature $\mathbb{F}_{\mathbf{I}}$ can be found by

$$\theta_{\mathbf{t}}^{opt} = \underset{\theta_{\mathbf{t}}}{\operatorname{argmax}} \mathcal{S}(\mathbb{F}_{\mathbf{I}}, \mathcal{R}(\mathbb{F}_{\mathbf{t}}, \theta_{\mathbf{t}})) \quad (6)$$

Substituting into Eq.5, we attain a simplified measurement model obviating θ and conditioned solely on \mathbf{t} as

$$P(\mathbf{I} | \mathbf{t}) = A' \cdot \mathcal{S}(\mathbb{F}_{\mathbf{I}}, \mathcal{R}(\mathbb{F}_{\mathbf{t}}, \theta_{\mathbf{t}}^{opt})) \quad (7)$$

For solving Eq.6, we rotate $\mathbb{F}_{\mathbf{t}}$ with uniformly sampled $\theta_{\mathbf{t}}$ in $[0, 2\pi)$, and keep the best. This discretized search initializes rotation to a rough value, which will later be refined as in §3.4. The rotation matching process is highly efficient since it reuses the same circular feature and does not render new hypotheses, where its throughput is detailed in Table.3.

3.2. Map Branch

In this section, we show how circular features are rendered from 2D floor maps with a given camera pose.

Rasterized 2D Point Cloud Map. Given a general 2D map representation \mathbf{M} (e.g. floorplan or occupancy grid) encoding area occupancy info, we uniformly sample points on the occupancy boundaries (i.e. walls) to form a 2D point cloud $\mathbb{M} = \{\mathbf{m}_i | i = 0 \cdots N-1\}$. Each point $\mathbf{m}_i = [\mathbf{t}_i, \mathbf{n}_i, \mathbf{s}_i]$ encodes its location \mathbf{t}_i , normal vector \mathbf{n}_i and optional semantic information \mathbf{s}_i . Available semantic information (e.g. door or window labels), are encoded as multiple binary masks appended to the point representation.

Latent Space Rendering. To circumvent the inefficient two-stage rendering-and-encoding process, we propose latent space rendering that directly renders circular features for given locations by aggregating features from visible map points. However, visibility proved to be a necessary but insufficient cue for the effective selection and rendering of latent space features (as shown in Fig.4). More specifically, visibility for static environments is locally constant at most sampling locations, providing a limited spatial context. To mitigate potential homogenization of our representations, we analyze fine-grain rendering dynamics such as length and incident-angle of the viewing rays between features and sampling location, and define an adaptive rendering mechanism.

Rendering Codebooks. We propose an over-specified latent space in order to endow map points with view-dependent features to encode rendering dynamics. We encode the 2D point cloud map with a 2D variant of PointNet

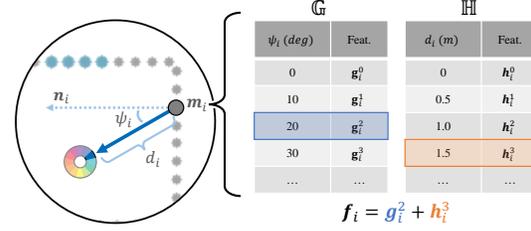


Figure 3. **Rendering codebook.** The map point features are dynamically (view-dependent) determined at rendering-time by the codebook and rendering dynamics. Linear interpolation is applied between adjacent codes in the codebook. Features from multiple codebooks are aggregated by summation.

[22] to assign each map point \mathbf{m}_i two sets of features $\mathbb{G}_i = \{\mathbf{g}_i^\beta | \beta = 0 \cdots G-1\}$ and $\mathbb{H}_i = \{\mathbf{h}_i^\gamma | \gamma = 0 \cdots H-1\}$, denoted as distance and incident-angle codebook respectively. Features in the codebook have the same dimension as circular feature segments $\mathbf{g}^\beta, \mathbf{h}^\gamma \in \mathbb{R}^D$. At rendering-time, map point features are chosen from the codebooks based on their distance and incident-angle w.r.t. the rendering location as shown in Fig.3. Formally, assume a rendering location $\hat{\mathbf{t}}$ and a map point $\mathbf{m}_i = [\mathbf{t}_i, \mathbf{n}_i, \mathbf{s}_i]$, let $\mathbf{d}_i = \mathbf{t}_i - \hat{\mathbf{t}}$, we can compute its rendering dynamics by

$$d_i = \|\mathbf{d}_i\| \quad (8)$$

$$\psi_i = \operatorname{atan2}(\|\mathbf{d}_i \times \mathbf{n}_i\|, \mathbf{d}_i \cdot \mathbf{n}_i) \quad (9)$$

where d_i and ψ_i are distance and incident-angle respectively. The clockwise incident-angle $\psi \in [0, 2\pi)$ distinguishes the four quadrants. With \mathbf{m}_i 's associated codebooks \mathbb{G}_i and \mathbb{H}_i , its feature \mathbf{f}_i is then determined by

$$\mathbf{f}_i = \mathbf{g}_i^{\frac{G\psi_i}{2\pi}} + \mathbf{h}_i^{\min(\frac{Hd_i}{d_{max}}, H)} \quad (10)$$

where d_{max} is a pre-defined maximum distance for the distance codebook. Similar to Eq.4, for non-integer indexing, we linearly interpolate between its two closest codes. Finally, if \mathbf{m}_i passes the visibility test to location $\hat{\mathbf{t}}$, we project \mathbf{f}_i to circular feature $\mathbb{F}_{\hat{\mathbf{t}}} = \{\mathbf{f}_{\hat{\mathbf{t}}}^\alpha | \alpha = 0 \cdots V-1\}$ by

$$\mathbf{f}_{\hat{\mathbf{t}}}^{\frac{V\omega_i}{2\pi}} = \mathbf{f}_i \quad (11)$$

where ω_i is the angle of viewing ray \mathbf{d}_i . Finally, the projected map point features are averaged into each segment. See supplementary for more rendering details.

3.3. Image Branch

In this section, we show how circular features are extracted from panorama and perspective images.

Circular Feature from Panoramas. For panorama images in equirectangular projection, each image column corresponds to a fixed horizontal FoV, as shown in Fig.2(b). Such

capture configuration facilitates a bijective mapping between our groups of adjacent input image columns and the segments in our rendered circular representation. The query panorama is fed into a ResNet50 [10] encoder to obtain a feature map, which is subsequently squeezed by averaged-pooling in the vertical dimension to comply with the feature dimensions of our circular segments, and averaged-pooled again in the horizontal direction to V elements, in accordance to the preconfigured number of feature segments in each circular feature.

Circular Feature from Perspective Images. We assume input perspective query images having known FoV and zero pitch/roll angle w.r.t. the ground plane. Note that for indoor query images, pitch/roll angles may be rectified from vanishing point estimates [6, 9]. Accordingly, each perspective image column corresponds to a non-fixed but known horizontal FoV as shown in Fig.2(a). We extract an image feature map with a ResNet50 encoder, using average pooling to squeeze the vertical dimension, and apply a perspective-to-equirectangular transform on the feature map to get the final circular feature. Since perspective images have no more than 180° FoV, its circular feature will have segments without assigned values, which will be masked out in the computation. Eq.3 will also be re-normalized to have range $[0, 1]$. In supplementary, our model’s robustness to capture pitch/roll angle alignment noise is detailed.

3.4. Refinement Branch

We address the discretized pose sampling nature of our MCL approach, by proposing a light-weight continuous refinement branch to improve upon the current estimation. As Fig.2(d) shows, with current best estimation \mathbf{t}^* and θ^* , our refinement branch takes two circular features \mathbb{F}_I and $\mathcal{R}(\mathbb{F}_{\mathbf{t}^*}, \theta^*)$ as input. The refinement network uses two 1D convolution layers with circular padding followed by a fully-connected layer to predict two offsets $\delta\mathbf{t}$, $\delta\theta$ for translation and rotation respectively. Then we render the updated map circular feature $\mathcal{R}(\mathbb{F}_{\mathbf{t}^*+\delta\mathbf{t}}, \theta^*+\delta\theta)$ and compute its similarity to \mathbb{F}_I using Eq.3. If the similarity score improved upon the original camera pose, we accept the step and iterate, otherwise we consider the refinement converged. The first refinement is always accepted to unquantize the estimation. This refinement usually converges within 3 iterations.

3.5. Training & Inference

Triplet Loss. We use triplet loss [27] to learn a mutual metric space between images and maps. To form a triplet, we let the image circular feature \mathbb{F}_I be the anchor, the map circular feature at ground truth camera pose $\mathbb{F}^+ = \mathcal{R}(\mathbb{F}_{\mathbf{t}_{gt}}, \theta_{gt})$ be the positive, and the map circular feature at a randomly sampled camera pose $\mathbb{F}^- = \mathcal{R}(\mathbb{F}_{\mathbf{t}_{rnd}}, \theta_{rnd})$ to serve as the negative. Then the triplet loss is defined as

$$\mathcal{L}_{triplet} = 2 \cdot \max(\mathcal{S}(\mathbb{F}_I, \mathbb{F}^+) - \mathcal{S}(\mathbb{F}_I, \mathbb{F}^-) + 0.5, 0) \quad (12)$$

Context Loss. Our similarity function \mathcal{S} , and consequently also our triplet loss $\mathcal{L}_{triplet}$, relies on aggregating element-wise comparisons, which effectively disregarding any intra-feature context. We design an additional context loss, providing feature segments a wider scope of its circular feature for learning context information (i.e. properties of the room/map). We first define circular feature context $\bar{\mathbb{F}}$ as the mean of its normalized feature segments

$$\bar{\mathbb{F}} = \frac{\sum_{\alpha=1}^V \mathbf{f}^\alpha / \|\mathbf{f}^\alpha\|}{V} \quad (13)$$

which we apply to our training triplet similarly to Eq.12

$$\mathcal{L}_{context} = \max(\cos(\bar{\mathbb{F}}_I, \bar{\mathbb{F}}^+) - \cos(\bar{\mathbb{F}}_I, \bar{\mathbb{F}}^-) + 1.0, 0) \quad (14)$$

With the context loss, circular features achieve better coarse level expressiveness, improving recall for query images with limited FoV, as shown in Table.2. This loss also acts as a regularizer by mitigating feature segments having large variance, leading to a smoother posterior estimation as shown in Fig.4.

Refinement Loss. For training the refinement branch, we sample circular features within a 0.5 meter radius and a 30 degree angle from the ground truth camera pose. We supervise the refinement branch using a regression loss as

$$\begin{aligned} \mathcal{L}_{refine.t} &= \|(\mathbf{t}_{gt} - \mathbf{t}^*) - \delta\mathbf{t}\| \\ \mathcal{L}_{refine.r} &= \min(|(\theta_{gt} - \theta^*) - \delta\theta|, 2\pi - |(\theta_{gt} - \theta^*) - \delta\theta|) \end{aligned} \quad (15)$$

Implementation Details. For triplet and context loss we sample 100 negative samples and broadcast the single ground truth (GT) sample for each training iteration. For refinement loss, we sample 20 hard negatives near the GT camera pose with a disturbance sampled from uniform distribution bounded in 30 degree and 0.5 meter radius. We combine the mean of all losses with equal weights. We set the hyper-parameters as $G = H = 32$, $V = 16$, $D = 128$ and $d_{max} = 10m$ consistently throughout the benchmarking. We sample the map into 2D point cloud with a 10 cm interval at occupancy boundaries. We render circular features for a $0.1m \times 0.1m$ uniform grid within the map range. For solving the relative rotations in Eq.6, we evaluate 16 uniformly sampled angles and keep the best. Finally, the posterior distribution is estimated using Eq.1,7. To extract final estimations from the posterior grid map, we apply a 3×3 non-maximum suppression to extract the maximums. For maximums that have larger score than a threshold (i.e. 0.8), we send them into the refinement branch to get the final estimations with their likelihoods as uncertainty estimation. Sorting by their likelihoods, a top-k estimation is available. More details see supplementary.

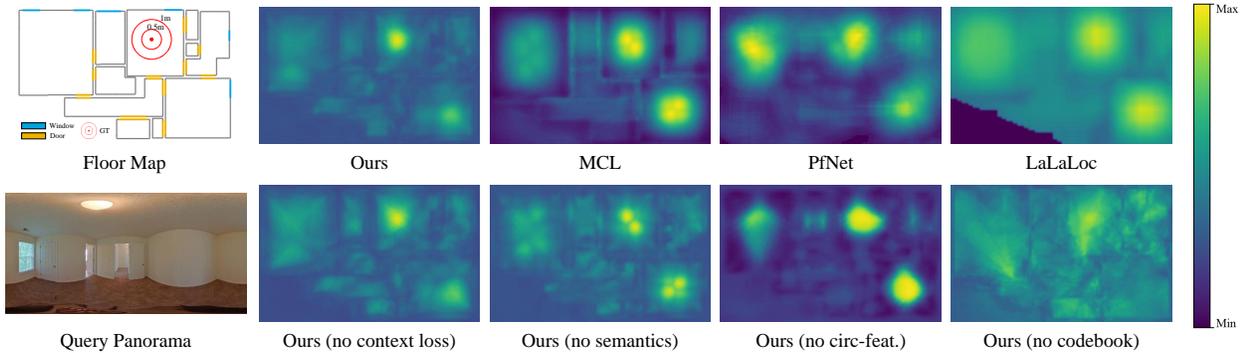


Figure 4. **Posterior map comparison with panorama query.** Our method leverages semantic information and provides a single clear maximum without ambiguities compared to baselines. Without context loss, the posterior map becomes slightly noisier. Without semantic labels, symmetric ambiguities emerge. Without circular feature, an accurate clear maximum cannot be identified from the clusters. Without codebook, the posterior map fails to show clear maxima.

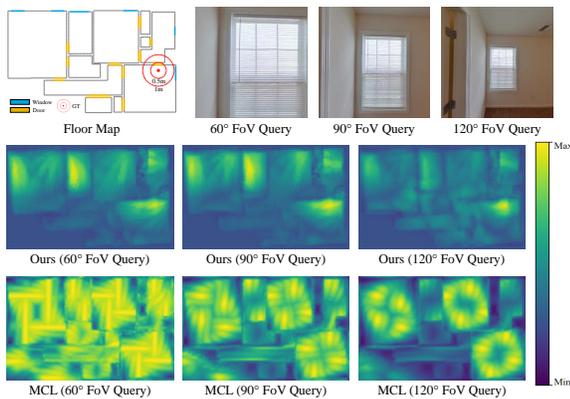


Figure 5. **Posterior map comparison with perspective queries.** With increasing query FoV, our method has increasing confidence around the GT location. Compared to MCL, our method leverages semantic information, which greatly reduces ambiguities.

4. Experiments

4.1. Datasets and Setups

We test on Zillow Indoor Dataset (ZInD) [5] and Structured3D dataset (S3D) [38]. ZInD provides 1,575 real world unfurnished residential homes with 59,361 panoramas, while S3D is a synthetic dataset containing 3,500 houses created by designers with 21,835 panoramas each having different lighting and furnishing levels. Both datasets provide 2D floor maps with windows and doors labels, and the 360° panoramas are registered to the floor maps.

We follow the official training/testing splitting for both datasets, and train separate models for panorama and perspective images. For benchmarking (Table.1 and Fig.7a), where we report performance for multiple perspective image FoVs, the model for perspective images is trained using images with FoVs uniformly sampled from 45° to 135°. For ablation study experiments, we report performance on

a representative 90° FoV, where the model is trained exclusively on images with 90° FoV. All perspective images are cropped from panoramas in the dataset with zero pitch, roll angles, a random yaw angle and equal horizontal/vertical FoVs.

4.2. Benchmark

We use LaLaLoc [12], PfNet [14] and MCL [7] as our baseline methods (see Table.1). For MCL, we simulate a 72-ray 2D LiDAR giving ground-truth distance without noise as its input. The PfNet takes input a semantic labelled floor map same as the one in Fig.4. For LaLaLoc, we follow their original protocol and take panoramas with known rotation as input and samples at $0.5m \times 0.5m$ grid. For PfNet and MCL, we sample a $0.1m \times 0.1m$ grid with 16 rotation angles same to our framework. We also attached a random baseline to show the statistics of the dataset. To compare with methods not utilizing semantic map labels (i.e. MCL, LaLaLoc), we report our performance without using semantic information as well. LASER exhibits superior accuracy and recall compared to our baseline methods. When the semantic labeling is not available, LASER recall rate is close to MCL with the GT depths as input. For panorama queries, LASER exhibits slightly better accuracy on S3D compared to ZInD due to the slight annotation errors in ZInD. For perspective query images, distance information becomes harder to extract since the room layout is sometimes not observable. In such case, LASER works better with furnished rooms in S3D, where the furniture provide cues for measuring distance. See supplementary for a more detailed discussion.

4.3. Performance Studies

Qualitative Study. In Figs.4,5, we visualize posterior maps for LASER and our baselines. Likelihoods are normalized w.r.t. their upper and lower bounds which vary among methods. We square our cosine distance for better visu-

| Query | Method | ZInD | | | | | | | Structured3D (Furnishing-Level : Full) | | | | | | |
|----------------------|---------------|-------------------|--------------------|-----------------|-----------------|---------------|---------------------|---------------------|--|--------------------|-----------------|-----------------|---------------|---------------------|---------------------|
| | | <1m med terr (cm) | <1m med rerr (deg) | 10cm recall (%) | 50cm recall (%) | 1m recall (%) | 1m & 30° recall (%) | top-3 1m recall (%) | <1m med terr (cm) | <1m med rerr (deg) | 10cm recall (%) | 50cm recall (%) | 1m recall (%) | 1m & 30° recall (%) | top-3 1m recall (%) |
| - | Random | (70.71) | (90.00) | 0.00 | 0.61 | 2.15 | 0.26 | 5.71 | (70.71) | (90.00) | 0.00 | 0.53 | 2.36 | 0.29 | 7.48 |
| Panorama | PfNet | 48.77 | 15.20 | 0.70 | 19.21 | 37.15 | 28.82 | 50.58 | 44.37 | 14.97 | 1.65 | 27.52 | 47.38 | 36.48 | 64.05 |
| | LaLaLoc | 10.65 | - | 35.61 | 71.69 | 76.00 | - | 91.62 | 6.83 | - | 58.57 | 85.98 | 87.51 | - | 98.23 |
| | MCL | 11.88 | 5.60 | 38.96 | 86.33 | 90.15 | 85.21 | 98.66 | 6.44 | 7.18 | 57.22 | 77.49 | 86.51 | 67.12 | 99.41 |
| | Ours (no sem) | 5.66 | 0.49 | 67.52 | 86.81 | 88.48 | 85.24 | 96.17 | 4.83 | 0.28 | 59.99 | 75.19 | 83.50 | 67.00 | 97.17 |
| | Ours | 5.16 | 0.47 | 78.83 | 96.83 | 97.12 | 96.99 | 98.90 | 3.87 | 0.23 | 79.20 | 95.05 | 95.52 | 94.76 | 98.41 |
| Perspective 60° FoV | PfNet | 63.31 | 23.48 | 0.21 | 5.23 | 15.86 | 9.19 | 24.87 | 61.27 | 17.60 | 0.35 | 6.01 | 16.91 | 11.20 | 26.75 |
| | MCL | 21.72 | 5.40 | 6.05 | 18.66 | 23.58 | 20.13 | 39.76 | 21.79 | 7.50 | 5.48 | 14.50 | 18.80 | 12.79 | 35.12 |
| | Ours (no sem) | 26.22 | 1.27 | 3.56 | 19.40 | 24.94 | 20.95 | 42.46 | 32.95 | 1.53 | 1.89 | 16.97 | 23.69 | 14.67 | 46.38 |
| | Ours | 29.39 | 1.21 | 4.61 | 34.42 | 44.30 | 42.10 | 61.88 | 16.97 | 0.98 | 11.02 | 41.48 | 45.96 | 43.13 | 65.65 |
| Perspective 90° FoV | PfNet | 56.55 | 14.24 | 0.55 | 10.97 | 26.78 | 19.63 | 38.62 | 60.00 | 12.49 | 0.59 | 8.78 | 24.34 | 18.74 | 34.41 |
| | MCL | 17.00 | 5.08 | 11.28 | 30.18 | 34.72 | 31.41 | 53.43 | 15.41 | 6.14 | 9.55 | 22.63 | 26.81 | 20.33 | 46.91 |
| | Ours (no sem) | 19.29 | 0.94 | 7.99 | 29.87 | 35.15 | 31.57 | 54.37 | 25.11 | 0.87 | 3.83 | 22.04 | 27.58 | 18.74 | 51.27 |
| | Ours | 22.09 | 0.85 | 9.42 | 51.62 | 59.40 | 57.77 | 76.23 | 12.99 | 0.64 | 20.86 | 53.51 | 56.28 | 54.21 | 76.31 |
| Perspective 120° FoV | PfNet | 53.81 | 11.74 | 0.74 | 14.21 | 31.42 | 25.94 | 43.83 | 60.27 | 12.51 | 0.77 | 10.14 | 28.05 | 22.39 | 38.95 |
| | MCL | 14.72 | 5.17 | 17.12 | 43.19 | 48.17 | 43.96 | 66.52 | 12.66 | 6.61 | 16.44 | 34.00 | 39.13 | 29.40 | 59.46 |
| | Ours (no sem) | 15.56 | 0.92 | 14.19 | 42.49 | 46.99 | 43.75 | 65.86 | 22.96 | 1.03 | 7.72 | 32.76 | 40.19 | 27.64 | 63.64 |
| | Ours | 19.07 | 0.80 | 15.18 | 65.37 | 72.14 | 71.08 | 85.28 | 11.31 | 0.70 | 30.88 | 68.83 | 70.95 | 69.77 | 87.98 |

Table 1. **Comparison with baselines on ZInD and S3D (fully furnished) dataset.** For indicating accuracy, we report median translation (terr) and rotation error (rerr) for all instances localized under 1m. We report recall at different translation accuracy levels, recall for inliers (<1m and <30°) and top-3 recall at 1m.

| Model | Panorama | | | Persp. 90° FoV | | |
|--------------|-------------------|--------------------|---------------|-------------------|--------------------|---------------|
| | <1m med terr (cm) | <1m med rerr (deg) | 1m recall (%) | <1m med terr (cm) | <1m med rerr (deg) | 1m recall (%) |
| base | 7.22 | 5.45 | 97.06 | 20.59 | 5.46 | 54.55 |
| + refine | 5.16 | 0.47 | 97.12 | 20.76 | 1.15 | 54.53 |
| - cxtloss | 8.54 | 5.45 | 95.76 | 21.57 | 5.46 | 50.52 |
| - pointnet | 7.84 | 5.46 | 95.87 | 23.46 | 5.68 | 50.26 |
| - codebook | 22.74 | 5.87 | 64.26 | 61.05 | 23.70 | 10.34 |
| - circ-feat. | 17.30 | - | 58.24 | 46.31 | - | 20.11 |

‘+’ with component ‘-’ without component

Table 2. **Ablation study over model components.** The base model uses the estimation from the posterior map without refinement. Without codebook, each map point is assigned a fixed feature. Without PointNet, all map points with the same semantic label share a fixed codebook. Without the circular structure (i.e. $V = 1$), the feature becomes agnostic to rotation.

alization. Fig.6 shows LASER is robust to challenging cases such as complex appearance and geometry conditions. While most failure cases are caused by ambiguities, some long-tail distributed location and texture also causes failure. **Model Ablation.** In Table.2, we show the ablation study over model components on ZInD dataset. Recall and median errors are reported for panorama and persp-90° input. Our refinement step unquantizes the discrete estimations, improving rotation accuracy from quantized initialization. Translation accuracy is similarly improved when the sampling density becomes a bottleneck. Context loss improves recall, and the improvements become more evident for queries with small FoV. Replacing the PointNet with a shared codebook across map points with same semantic labels makes LASER agnostic to the input map domain, but slightly degrades all metrics. Omission of codebook rendering or circular feature encoding, largely degrades all performance metrics.

Performance over FoV. In Fig.7a, we show performance over different query image FoVs. The perspective query images have same horizontal and vertical FoV, while equirectangular query images always have 180° vertical FoV. With

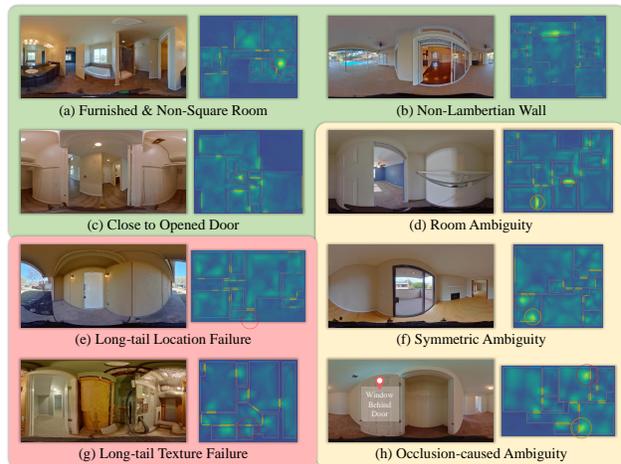


Figure 6. **Qualitative study on method robustness under challenging cases.** Success, ambiguous and failure cases are placed inside green, yellow and red boxes respectively. The GT locations are circled red and ambiguities are circled yellow in the posterior maps where floor maps are overlaid.

increasing image FoV, LASER consistently gains better performance. LASER has lower translation error with equirectangular views, but lower rotation error from perspective views, while their 1m-recall is similar.

Performance over Hyper-Parameters. Fig.7b shows that the incident-angle codebook improves rotation accuracy while the distance codebook is better in translation. Combining both gives the best. Performance is not sensitive to codebook size, where a modest number (i.e. 16) is satisfactory. Fig.7c shows increasing the translation sample density improves the recall and translation accuracy, but such improvement becomes marginal after $0.1m \times 0.1m$. While recall is not sensitive to rotation sample density, where a modest number (i.e. 16) is satisfactory. Fig.7d shows a modest rendering resolution V (i.e. $16/32$) is satisfactory.

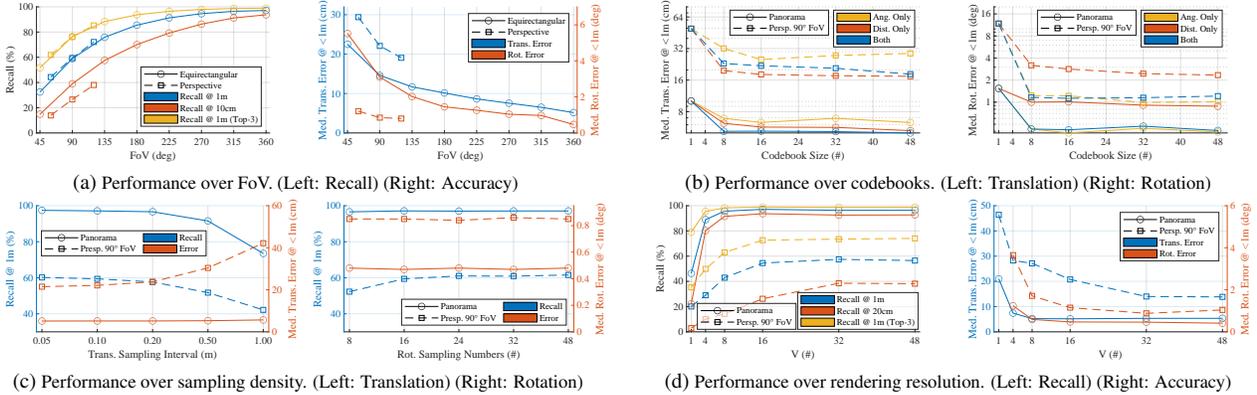


Figure 7. **Detailed performance analysis.** Different query types are shown in line-types. Different metrics/configurations are shown in colors. Note that some figures have two y-axes showing different metrics in separate colors.

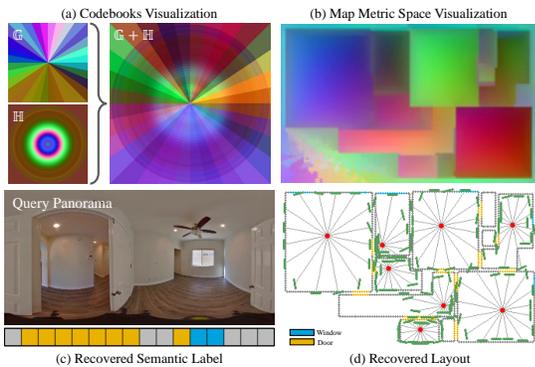


Figure 8. **Visualization and interpretability.** (a) Codebook visualization using cosine kernel PCA. (b) Visualization of map metric space using cosine kernel t-SNE. (c,d) Recovered image semantic labels and room layout by inversely matching image circular feature exhaustively against all map point codebooks. (d) The recovered viewing rays and incident-angles are shown in gray and green lines respectively.

Codebook Visualization and Interpretability. As shown in Fig.8(a,b), both codebook and sampled feature map exhibit smooth transitions between adjacent codes while there are distinctiveness for distant codes and different rooms. The distance codebook has less variation after a certain distance, which is a cue for choosing d_{max} . Furthermore, we exhaustively match image feature to the codebooks of all map points. We greedily record the code with closest distance, where the semantic label can be recovered from its map point as shown in Fig.8(c). The distance and incident-angle can be also recovered w.r.t. to its codebook as shown in Fig.8(d). This shows our model is implicitly learning to extract semantic and layout from the image.

Timing. As shown in Table.3, LASER is significantly faster than existing works both in sampling and query, which allows its application in time-sensitive re-localization. Within the query, the refinement step is relatively slower since the sequential rendering is not massively parallelizable.

| Method | Sampling Fps | Sampling Time (s) | Query Fps | Query Module | Time (ms) |
|---------|--------------|-------------------|-------------|--------------|-------------|
| PfNet | 2471 ± 1328 | 48.95 ± 38.95 | 5.06 ± 1.77 | ResNet50 | 10.5 ± 1.2 |
| LaLaLoc | 15.47 ± 5.96 | 25.13 ± 23.67 | 0.30 ± 0.06 | Measure | 17.7 ± 7.4 |
| Ours | 13238 ± 1890 | 0.97 ± 1.09 | 8.31 ± 0.64 | Refine | 92.7 ± 16.1 |

Table 3. **Timing.** Performance on S3D with a single NVIDIA Tesla V100 following the algorithm configurations as in §4.2.

5. Conclusions and Discussions

LASER introduces the concept of feature codebook which leverages the redundancy of overspecified feature predictions to achieve runtime dynamic latent feature assignments, while obviating additional expensive encoding processes. Moreover, we instantiate our codebook scheme into a novel latent space rendering process, where rendering dynamics are efficiently encoded into latent feature representations. In practice, latent space rendering enables our proposed geometrically-structured metric learning framework to achieve state-of-art efficiency and accuracy within the 2D visual localization task.

Besides the success that feature codebook achieved in the 2D visual localization task, the codebook scheme can also be applied in general learning tasks to replace the encoder network. In addition to the speed-up, the codebook scheme enables the seamless debugging using the (nearest-neighbor) inverse matching as shown in Fig.8(c,d). In practice, the nearest-neighbor matching can be easily extended to extract probabilistic estimation. We believe such tool has strong potential to benefit general learning tasks in the aspects of interpretability, probabilistic estimation/calibration, and improving speed/latency.

Limitations to address in future work include modeling object-based (e.g. furniture) partial map occlusions. Moreover, we have not considered the semantics between visibility and environmental state (e.g. leveraging a door’s open/close status to reason about the content visibility). Finally, the scope of our latent space rendering, may be extended to 3D space or consider more complex rendering-time dynamics such as surface BRDF.

References

- [1] Vassileios Balntas, Shuda Li, and Victor Prisacariu. Relocnet: Continuous metric learning relocalisation using neural nets. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 751–767, 2018. 2
- [2] Federico Boniardi, Abhinav Valada, Rohit Mohan, Tim Caselitz, and Wolfram Burgard. Robot localization in floor plans using a room layout edge extraction network. In *2019 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 5291–5297. IEEE, 2019. 1, 2
- [3] Eric Brachmann, Alexander Krull, Sebastian Nowozin, Jamie Shotton, Frank Michel, Stefan Gumhold, and Carsten Rother. Dsac-differentiable ransac for camera localization. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 6684–6692, 2017. 2
- [4] Hang Chu, Dong Ki Kim, and Tsuhan Chen. You are here: Mimicking the human thinking process in reading floor-plans. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 2210–2218, 2015. 2
- [5] Steve Cruz, Will Hutchcroft, Yuguang Li, Naji Khosravan, Ivaylo Boyadzhiev, and Sing Bing Kang. Zillow indoor dataset: Annotated floor plans with 360deg panoramas and 3d room layouts. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2133–2143, 2021. 1, 2, 6
- [6] Benjamin Davidson, Mohsan S Alvi, and João F Henriques. 360° camera alignment via segmentation. In *Computer Vision—ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XXVIII 16*, pages 579–595. Springer, 2020. 5
- [7] Frank Dellaert, Dieter Fox, Wolfram Burgard, and Sebastian Thrun. Monte carlo localization for mobile robots. In *Proceedings 1999 IEEE International Conference on Robotics and Automation (Cat. No. 99CH36288C)*, volume 2, pages 1322–1328. IEEE, 1999. 1, 2, 3, 6
- [8] Mingyu Ding, Zhe Wang, Jiankai Sun, Jianping Shi, and Ping Luo. Camnet: Coarse-to-fine retrieval for camera relocalization. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 2871–2880, 2019. 2
- [9] Erwan Guillou, Daniel Meneveau, Eric Maisel, and Kadi Bouatouch. Using vanishing points for camera calibration and coarse 3d reconstruction from a single image. *The Visual Computer*, 16(7):396–410, 2000. 5
- [10] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016. 5
- [11] Elad Hoffer and Nir Ailon. Deep metric learning using triplet network. In *International workshop on similarity-based pattern recognition*, pages 84–92. Springer, 2015. 3
- [12] Henry Howard-Jenkins, Jose-Raul Ruiz-Sarmiento, and Victor Adrian Prisacariu. Lalaloc: Latent layout localisation in dynamic, unvisited environments. *arXiv preprint arXiv:2104.09169*, 2021. 1, 2, 6
- [13] Rico Jonschkowski, Divyam Rastogi, and Oliver Brock. Differentiable particle filters: End-to-end learning with algorithmic priors. *arXiv preprint arXiv:1805.11122*, 2018. 2
- [14] Peter Karkus, David Hsu, and Wee Sun Lee. Particle filter networks with application to visual localization. In *Conference on robot learning*, pages 169–178. PMLR, 2018. 1, 2, 6
- [15] Alex Kendall, Matthew Grimes, and Roberto Cipolla. Posenet: A convolutional network for real-time 6-dof camera relocalization. In *Proceedings of the IEEE international conference on computer vision*, pages 2938–2946, 2015. 2
- [16] Ricardo Martin-Brualla, Noha Radwan, Mehdi SM Sajjadi, Jonathan T Barron, Alexey Dosovitskiy, and Daniel Duckworth. Nerf in the wild: Neural radiance fields for unconstrained photo collections. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7210–7219, 2021. 2
- [17] Oscar Mendez, Simon Hadfield, Nicolas Pugeault, and Richard Bowden. Sedar-semantic detection and ranging: Humans can localise without lidar, can robots? In *2018 IEEE International Conference on Robotics and Automation (ICRA)*, pages 6053–6060. IEEE, 2018. 2
- [18] 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 *European conference on computer vision*, pages 405–421. Springer, 2020. 2
- [19] Zhixiang Min and Enrique Dunn. Voldor+ slam: For the times when feature-based or direct methods are not good enough. In *2021 IEEE International Conference on Robotics and Automation (ICRA)*, pages 13813–13819. IEEE, 2021. 2
- [20] Zhixiang Min, Yiding Yang, and Enrique Dunn. Voldor: Visual odometry from log-logistic dense optical flow residuals. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4898–4909, 2020. 2
- [21] Jaesik Park, Qian-Yi Zhou, and Vladlen Koltun. Colored point cloud registration revisited. In *Proceedings of the IEEE international conference on computer vision*, pages 143–152, 2017. 1
- [22] Charles R Qi, Hao Su, Kaichun Mo, and Leonidas J Guibas. Pointnet: Deep learning on point sets for 3d classification and segmentation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 652–660, 2017. 2, 4
- [23] Paul-Edouard Sarlin, Cesar Cadena, Roland Siegwart, and Marcin Dymczyk. From coarse to fine: Robust hierarchical localization at large scale. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12716–12725, 2019. 2
- [24] Paul-Edouard Sarlin, Daniel DeTone, Tomasz Malisiewicz, and Andrew Rabinovich. Superglue: Learning feature matching with graph neural networks. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 4938–4947, 2020. 2
- [25] Paul-Edouard Sarlin, Ajaykumar Unagar, Mans Larsson, Hugo Germain, Carl Toft, Viktor Larsson, Marc Pollefeys, Vincent Lepetit, Lars Hammarstrand, Fredrik Kahl, et al. Back to the feature: Learning robust camera localization from pixels to pose. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3247–3257, 2021. 2

- [26] Johannes L Schonberger and Jan-Michael Frahm. Structure-from-motion revisited. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4104–4113, 2016. 1, 2
- [27] Florian Schroff, Dmitry Kalenichenko, and James Philbin. Facenet: A unified embedding for face recognition and clustering. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 815–823, 2015. 5
- [28] Jiaming Sun, Zehong Shen, Yuang Wang, Hujun Bao, and Xiaowei Zhou. Loftr: Detector-free local feature matching with transformers. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8922–8931, 2021. 2
- [29] Hajime Taira, Masatoshi Okutomi, Torsten Sattler, Mircea Cimpoi, Marc Pollefeys, Josef Sivic, Tomas Pajdla, and Akihiko Torii. Inloc: Indoor visual localization with dense matching and view synthesis. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 7199–7209, 2018. 2
- [30] Chengzhou Tang and Ping Tan. Ba-net: Dense bundle adjustment network. *arXiv preprint arXiv:1806.04807*, 2018. 2
- [31] Zachary Teed and Jia Deng. Deepv2d: Video to depth with differentiable structure from motion. *arXiv preprint arXiv:1812.04605*, 2018. 2
- [32] Shenlong Wang, Sanja Fidler, and Raquel Urtasun. Lost shopping! monocular localization in large indoor spaces. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 2695–2703, 2015. 1, 2
- [33] Xingkui Wei, Yinda Zhang, Zhuwen Li, Yanwei Fu, and Xiangyang Xue. DeepSfM: Structure from motion via deep bundle adjustment. In *European conference on computer vision*, pages 230–247. Springer, 2020. 2
- [34] Wera Winterhalter, Freya Fleckenstein, Bastian Steder, Luciano Spinello, and Wolfram Burgard. Accurate indoor localization for rgb-d smartphones and tablets given 2d floor plans. In *2015 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 3138–3143. IEEE, 2015. 1, 2
- [35] Luwei Yang, Ziqian Bai, Chengzhou Tang, Honghua Li, Yasutaka Furukawa, and Ping Tan. Sanet: Scene agnostic network for camera localization. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 42–51, 2019. 2
- [36] Kwang Moo Yi, Eduard Trulls, Yuki Ono, Vincent Lepetit, Mathieu Salzmann, and Pascal Fua. Learning to find good correspondences. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2666–2674, 2018. 2
- [37] Enliang Zheng and Changchang Wu. Structure from motion using structure-less resection. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 2075–2083, 2015. 2
- [38] Jia Zheng, Junfei Zhang, Jing Li, Rui Tang, Shenghua Gao, and Zihan Zhou. Structured3d: A large photo-realistic dataset for structured 3d modeling. In *Computer Vision—ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part IX 16*, pages 519–535. Springer, 2020. 1, 2, 6
- [39] Huizhong Zhou, Benjamin Ummenhofer, and Thomas Brox. Deeptam: Deep tracking and mapping. In *Proceedings of the European Conference on Computer Vision (ECCV)*, September 2018. 2
- [40] Qunjie Zhou, Torsten Sattler, and Laura Leal-Taixe. Patch2pix: Epipolar-guided pixel-level correspondences. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4669–4678, 2021. 2