

# 6D Camera Relocalization in Ambiguous Scenes via Continuous Multimodal Inference

Mai Bui<sup>1</sup> Tolga Birdal<sup>2</sup> Haowen Deng<sup>1,3</sup> Shadi Albarqouni<sup>1,4</sup>  
Leonidas Guibas<sup>2</sup> Slobodan Ilic<sup>1,3</sup> Nassir Navab<sup>1,5</sup>

<sup>1</sup> Technical University of Munich <sup>2</sup> Stanford University  
<sup>3</sup> Siemens AG <sup>4</sup> ETH Zurich <sup>5</sup> Johns Hopkins University

**Abstract.** We present a multimodal camera relocalization framework that captures ambiguities and uncertainties with continuous mixture models defined on the manifold of camera poses. In highly ambiguous environments, which can easily arise due to symmetries and repetitive structures in the scene, computing one plausible solution (what most state-of-the-art methods currently regress) may not be sufficient. Instead we predict multiple camera pose hypotheses as well as the respective uncertainty for each prediction. Towards this aim, we use Bingham distributions, to model the orientation of the camera pose, and a multivariate Gaussian to model the position, with an end-to-end deep neural network. By incorporating a Winner-Takes-All training scheme, we finally obtain a mixture model that is well suited for explaining ambiguities in the scene, yet does not suffer from mode collapse, a common problem with mixture density networks. We introduce a new dataset specifically designed to foster camera localization research in ambiguous environments and exhaustively evaluate our method on synthetic as well as real data on both ambiguous scenes and on non-ambiguous benchmark datasets. We plan to release our code and dataset under [multimodal3dvision.github.io](https://github.com/multimodal3dvision).

## 1 Introduction

Camera relocalization is the term for determining the 6-DoF rotation and translation parameters of a camera with respect to a known 3D world. It is now a key technology in enabling a multitude of applications such as augmented reality, autonomous driving, human computer interaction and robot guidance, thanks to its extensive integration in simultaneous localization and mapping (SLAM) [18, 25, 71], structure from motion (SfM) [75, 81], metrology [6] and visual localization [63, 76]. For decades, vision scholars have worked on finding the unique solution of this problem [41, 42, 42, 66, 73, 84, 85]. However, this trend is now witnessing a fundamental challenge. A recent school of thought has begun to point out that for highly complex and ambiguous real environments, obtaining a single solution for the location and orientation of a camera is simply not sufficient. This has led to a paradigm shift towards estimating a range of solutions, in the form of full probability distributions [1, 7, 8] or at least solutions that estimate the uncertainty in orientation estimates [43, 54]. Thanks to advances

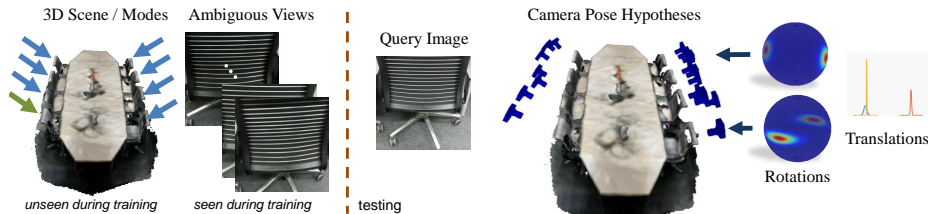


Fig. 1: In a highly ambiguous environment, similar looking views can easily confuse current camera pose regression models and lead to incorrect localization results. Instead, given a query RGB image, our aim is to predict the possible modes as well as the associated uncertainties, which we model by the parameters of Bingham and Gaussian mixture models.

in state-of-the-art machine learning, this important problem can now be tackled via data driven algorithms that are able to discover multi-modal and complex distributions, targeting the task at hand.

In this paper, we devise a multi-hypotheses method, depicted in Fig. 1, for learning continuous mixture models on manifold valued rotations (parameterized by quaternions) and Euclidean translations that can explain uncertainty and ambiguity in 6DoF camera relocalization, while avoiding mode collapse [72]. In particular, we propose a direct regression framework utilizing a combination of antipodally symmetric Bingham [4] and Gaussian probability distributions in order to deal with rotational and translational uncertainties respectively. Together, they are well suited to the geometric nature of  $SE(3)$  pose representations. Using said distributions, we aim to build novel variational models that enable estimation of full covariances on the discrete modes to be predicted. For better exploration of the variational posterior, we extend the established particle based training approaches [53, 53, 69] to mixtures of Gaussians and Bingham. While these techniques only optimize the location of individual hypothesis to cover diverse modes, we additionally learn to predict associated variances on the manifold. We can then approximate the manifold valued posterior of our problem in a continuous fashion. Note that, to the best of our knowledge, such continuous distributions for multi-modal modeling of the network posteriors formed by the 6D pose parameters have not been explored previously. Our synthetic and real experiments demonstrate promising performance both under ambiguities and non-ambiguous scenes. Our method is also flexible in the sense that it can be used with a wide variety of backbone architectures. In a nutshell, our contributions are:

1. We provide a general framework for continuously modelling conditional density functions on quaternions using Bingham distributions, while explaining the translational uncertainty with multi-modal Gaussians.
2. We propose a new training scheme for deep neural networks that enables the inference of a diverse set of modes and related concentration parameters as well as the prior weights for the mixture components.

3. We exhaustively evaluate our method on existing datasets, demonstrating the validity of our approach. Additionally, we create a new highly ambiguous camera relocalization dataset, which we use to showcase the quality of results attained by our algorithm and provide a thorough study on camera localization in ambiguous scenes.

## 2 Prior Art

6D camera relocalization is a very well studied topic with a vast literature [11, 13, 14, 16, 27, 46, 57, 76]. Our work considers the uncertainty aspect and this is what we focus on here: We first review the uncertainty estimation in deep networks, and subsequently move to uncertainty in 6D pose and relocalization.

**Characterizing the Posterior in Deep Networks** Typical CNNs [39, 77] are over-confident in their predictions [36, 88]. Moreover, these networks tend to approximate the conditional averages of the target data [9]. These undesired properties render the immediate outputs of those networks unsuitable for the quantification of calibrated uncertainty. This has fostered numerous works as we will summarize in the following. *Mixture Density Networks (MDN)* [9] is the pioneer to model the conditional distribution by predicting the parameters of a mixture model. Yet, it is repeatedly reported that optimizing for general mixture models suffers from mode collapse and numerical instabilities [22, 53]. These issues can to a certain extent be alleviated by using Dropout [30] as a Bayesian approximation, but even for moderate dimensions these methods still face difficulties in capturing multiple modes. Instead, the more tailored *Winner Takes All (WTA)* [28, 37] as well as *Relaxed-WTA (RWTA)* [69] try to capture the multimodal posterior in the  $K$ -best hypotheses predictions of the network. *Evolving-WTA (EWTA)* [53] further avoids the inconsistencies related to the WTA losses. Though, a majority of these works consider only low dimensional posterior with the assumption of a Euclidean space, whereas we consider a 7D non-Euclidean highly non-convex posterior.

**Uncertainty in 6D** Initial attempts to capture the uncertainty of camera relocalization involved Random Forests (RF) [15]. Valentin *et al.* [82] stored GMM components at the leaves of a scene coordinate regression forest [76]. The modes are obtained via a mean shift procedure, and the covariance is explained by a 3D Gaussian. A similar approach later considered the uncertainty in object coordinate labels [12]. It is a shortcoming of RF that both of these approaches require hand crafted depth features. Moreover, their uncertainty is on the correspondences and not on the final camera pose. Thus a costly RANSAC [29] is required to propagate the uncertainty in the leaves to the camera pose.

Only recently, Manhardt *et al.* [54] localized a camera against a known object under rotational ambiguities arising due to symmetries or self-occlusions. They extended the RTWA [69] to deal with the 3D rotations using quaternions. This method can only yield discrete hypotheses not continuous density estimates.

Similarly, the pose estimation network of Pitteri *et al.* [64] explicitly considered axis-symmetric objects whose pose cannot be uniquely determined. Likewise, Corona *et al.* [21] addressed general rotational symmetries. All of these works require extensive knowledge about the object and cannot be extended to the scenario of localizing against a scene without having a 3D model. Note that the latter two works cannot handle the case of self-symmetry and [21] additionally requires a dataset of symmetry-labeled objects, an assumption unlikely to be fulfilled in real applications.

Bayesian PoseNet [43] was one of the first works to model uncertainty for the 6D relocalization problem. It leveraged Dropout [30] to sample the posterior as a way to enable approximate probabilistic pose inference. Although in theory this method can generate discrete samples from the multi-modal distribution, in practice, as we will demonstrate, the Monte Carlo scheme tends to draw samples around a single mode. This method also suffers from the large errors associated to PoseNet [46] itself. The successive VidLoc [20] adapted MDNs [9] to model and predict uncertainty for the 6D relocalization problem. Besides the reported issues of MDNs, VidLoc incorrectly modeled the rotation parameters using Gaussians and lacked the demonstrations of uncertainty on rotations. Contrarily, in this work we devise a principled method using Bingham distributions [4] that are well suited to the double covering nature of unit quaternions. HydraNet [62] provided calibrated uncertainties on the  $SO(3)$  group, but assumed a unimodal posterior that is centered on the naive  $\mathbb{R}^4$ -mean of predicted quaternions.

Our work is inspired by [65], where a variational auto-encoder [48] is learnt to approximate the posterior of  $SO(2)$  modeled by von Mises mixtures [56]. Though, it is not trivial to tweak and generalize [65] to the continuous, highly multi-modal and multi-dimensional setting of 6D camera relocalization. This is what we precisely contribute in this work. Note that we are particularly interested in the *aleatoric* uncertainty (noise in the observations) and leave the *epistemic* (noise in the model) part as a future work [45].

### 3 The Bingham Distribution

Derived from a zero-mean Gaussian, the Bingham distribution [4] (BD) is an antipodally symmetric probability distribution conditioned to lie on  $\mathbb{S}^{d-1}$  with probability density function (PDF)  $\mathcal{B} : \mathbb{S}^{d-1} \rightarrow \mathbb{R}$ :

$$\mathcal{B}(\mathbf{x}; \mathbf{\Lambda}, \mathbf{V}) = (1/F) \exp(\mathbf{x}^T \mathbf{V} \mathbf{\Lambda} \mathbf{V}^T \mathbf{x}) = (1/F) \exp\left(\sum_{i=1}^d \lambda_i (\mathbf{v}_i^T \mathbf{x})^2\right) \quad (1)$$

where  $\mathbf{V} \in \mathbb{R}^{d \times d}$  is an orthogonal matrix ( $\mathbf{V} \mathbf{V}^T = \mathbf{V}^T \mathbf{V} = \mathbf{I}_{d \times d}$ ) describing the orientation,  $\mathbf{\Lambda} \in \mathbb{R}^{d \times d}$  is called the *concentration matrix* with  $0 \geq \lambda_1 \geq \dots \geq \lambda_{d-1}$ :  $\mathbf{\Lambda} = \text{diag}([0 \ \lambda_1 \ \lambda_2 \ \dots \ \lambda_{d-1}])$ .

It is easy to show that adding a multiple of the identity matrix  $\mathbf{I}_{d \times d}$  to  $\mathbf{V}$  does not change the distribution [4]. Thus, we conveniently force the first entry of  $\mathbf{\Lambda}$  to be zero. Moreover, since it is possible to swap columns of  $\mathbf{\Lambda}$ , we can build  $\mathbf{V}$  in a sorted fashion. This allows us to obtain *the mode* very easily by taking

the first column of  $\mathbf{V}$ . Due to its antipodally symmetric nature, the mean of the distribution is always zero.  $F$  in Eq (1) denotes the *normalization constant* dependent only on  $\mathbf{\Lambda}$  and is of the form:

$$F \triangleq |S_{d-1}| \cdot {}_1F_1\left(1/2, d/2, \mathbf{\Lambda}\right), \quad (2)$$

where  $|S_{d-1}|$  is the surface area of the  $d$ -sphere and  ${}_1F_1$  is a confluent hypergeometric function of matrix argument [40, 50]. The computation of  $F$  is not trivial. In practice, following Glover [34], this quantity as well as its gradients are approximated by tri-linear interpolation using a pre-computed look-up table over a predefined set of possible values in  $\mathbf{\Lambda}$ , lending itself to differentiation [49, 51].

**Relationship to quaternions** The antipodal symmetry of the PDF makes it amenable to explain the topology of quaternions, i. e.,  $\mathcal{B}(\mathbf{x}; \cdot) = \mathcal{B}(-\mathbf{x}; \cdot)$  holds for all  $\mathbf{x} \in \mathbb{S}^{d-1}$ . In 4D when  $\lambda_1 = \lambda_2 = \lambda_3$ , one can write  $\mathbf{\Lambda} = \text{diag}([1, 0, 0, 0])$ . In this case, Bingham density relates to the dot product of two quaternions  $\mathbf{q}_1 \in \mathbb{H}_1 \triangleq \mathbf{x}$  and the mode of the distribution, say  $\bar{\mathbf{q}}_2 \in \mathbb{H}_1$ . This induces a metric of the form:  $d_{\text{bingham}} = d(\mathbf{q}_1, \bar{\mathbf{q}}_2) = (\mathbf{q}_1 \cdot \bar{\mathbf{q}}_2)^2 = \cos(\theta/2)^2$ .

Bingham distributions have been extensively used to represent distributions on unit quaternions ( $\mathbb{H}_1$ ) [5, 8, 32, 33, 50]; however, to the best of our knowledge, never for the problem we consider here.

**Constructing a Bingham distribution on a given mode** Creating a Bingham distribution on any given mode  $\mathbf{q} \in \mathbb{H}_1$  requires finding a set of vectors orthonormal to  $\mathbf{q}$ . This is a frame bundle  $\mathbb{H}_1 \rightarrow \mathcal{F}\mathbb{H}_1$  composed of four unit vectors: the mode and its orthonormals. We follow Birdal *et al.* [8] and use the *parallelizability* of unit quaternions to define the orthonormal basis  $\mathbf{V} : \mathbb{H}_1 \mapsto \mathbb{R}^{4 \times 4}$ :

$$\mathbf{V}(\mathbf{q}) \triangleq \begin{bmatrix} q_1 & -q_2 & -q_3 & q_4 \\ q_2 & q_1 & q_4 & q_3 \\ q_3 & -q_4 & q_1 & -q_2 \\ q_4 & q_3 & -q_2 & -q_1 \end{bmatrix}. \quad (3)$$

It is easy to verify that the matrix valued function  $\mathbf{V}(\mathbf{q})$  is orthonormal for every  $\mathbf{q} \in \mathbb{H}_1$ .  $\mathbf{V}(\mathbf{q})$  further gives a convenient way to represent quaternions as matrices paving the way to linear operations, such as quaternion multiplication or orthonormalization without the Gram-Schmidt.

**Relationship to other representations** Note that geometric [3] or measure theoretic [26], there are multitudes of ways of defining probability distributions on the Lie group of 6D rigid transformations [38]. A choice would be to define Gaussian distribution on the Rodrigues vector (or exponential coordinates) [60] where the geodesics are straight lines [59] or the use of Concentrated Gaussian distributions [10] on matrices of  $\text{SE}(3)$ . However, as our purpose is not tracking but direct regression, in this work we favor quaternions as continuous and minimally redundant parameterizations without singularities [17, 35] and use the Bingham distribution that is well suited to their topology. We handle the redundancy  $\mathbf{q} \equiv -\mathbf{q}$  by mapping all the rotations to the northern hemisphere.

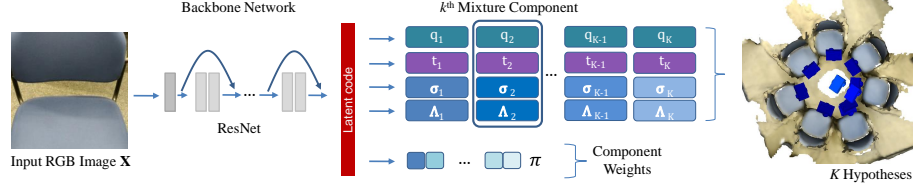


Fig. 2: Forward pass of our network. For an input RGB image we predict  $K$  camera pose hypotheses as well as Bingham concentration parameters, Gaussian variances and component weights to obtain a mixture model.

## 4 Proposed Model

We now describe our model for uncertainty prediction following [65]. We consider the situation where we observe an input image  $\mathbf{X} \in \mathbb{R}^{W \times H \times 3}$  and assume the availability of a predictor function  $\mu_{\mathbf{\Gamma}}(\mathbf{X}) : \mathbb{R}^{W \times H \times 3} \mapsto \mathbb{H}_1$  parameterized by  $\mathbf{\Gamma} = \{\mathbf{\Gamma}_i\}$ . Note that predicting entities that are non-Euclidean easily generalizes to prediction of Euclidean quantities such as translations e.g.  $\mathbf{t} \in \mathbb{R}^3$ . For the sake of conciseness and clarity, we will omit the translations and concentrate on the rotations. Translations modeled via Gaussians will be precised later on.

**The unimodal case** We momentarily assume that  $\mu_{\mathbf{\Gamma}}(\cdot)$ , or short  $\mu(\cdot)$ , can yield the correct values of the absolute camera rotations  $\mathbf{q}_i \in \mathbb{H}_1$  with respect to a common origin, admitting a non-ambiguous prediction, hence a posterior of single mode. We use the predicted rotation to set the most likely value (mode) of a BD:

$$p_{\mathbf{\Gamma}}(\mathbf{q} | \mathbf{X}; \mathbf{\Lambda}) = (1/F) \exp(\mathbf{q}^\top \mathbf{V}_\mu \mathbf{\Lambda} \mathbf{V}_\mu^\top \mathbf{q}), \quad (4)$$

and let  $\mathbf{q}_i$  differ from this value up to the extent determined by  $\mathbf{\Lambda} = \{\lambda_i\}$ . For the sake of brevity we use  $\mathbf{V}_\mu \equiv \mathbf{V}(\mu(\mathbf{X}))$ , the orthonormal basis aligned with the predicted quaternion  $\mu(\mathbf{X})$  and as defined in Eq (3).

While for certain applications, fixing  $\mathbf{\Lambda}$  can work, in order to capture the variation in the input, it is recommended to adapt  $\mathbf{\Lambda}$  [65]. Thus, we introduce it among the unknowns. To this end we define the function  $\mathbf{\Lambda}_{\mathbf{\Gamma}}(\mathbf{X})$  or in short  $\mathbf{\Lambda}_{\mathbf{\Gamma}}$  for computing the concentration values depending on the current image and the parameters  $\mathbf{\Gamma}$ . Our final model for the unimodal case reads:

$$p_{\mathbf{\Gamma}}(\mathbf{q} | \mathbf{X}) = \frac{\exp(\mathbf{q}^\top \mathbf{V}(\mu(\mathbf{X})) \mathbf{\Lambda}_{\mathbf{\Gamma}}(\mathbf{X}) \mathbf{V}(\mu(\mathbf{X}))^\top \mathbf{q})}{F(\mathbf{\Lambda}_{\mathbf{\Gamma}}(\mathbf{X}))} = \frac{\exp(\mathbf{q}^\top \mathbf{V}_\mu \mathbf{\Lambda}_{\mathbf{\Gamma}} \mathbf{V}_\mu^\top \mathbf{q})}{F(\mathbf{\Lambda}_{\mathbf{\Gamma}})} \quad (5)$$

The latter follows from the short-hand notations and is included for clarity. Given a collection of observations i.e., images  $\mathcal{X} = \{\mathbf{X}_i\}$  and associated rotations  $\mathbf{Q} = \{\mathbf{q}_i\}$ , where  $i = 1, \dots, N$ , the parameters of  $\mu_{\mathbf{\Gamma}}(\mathbf{X})$  and  $\mathbf{\Lambda}_{\mathbf{\Gamma}}(\mathbf{X})$  can be

obtained simply by maximizing the log-likelihood:

$$\mathbf{\Gamma}^* = \arg \max_{\mathbf{\Gamma}} \log \mathcal{L}_u(\mathbf{\Gamma}|\mathcal{X}, \mathbf{Q}) \quad (6)$$

$$\log \mathcal{L}_u(\mathbf{\Gamma}|\mathcal{X}, \mathbf{Q}) = \sum_{i=1}^N \mathbf{q}_i^\top \mathbf{V}_\mu \mathbf{\Lambda}_\mathbf{\Gamma} \mathbf{V}_\mu^\top \mathbf{q}_i - \sum_{i=1}^N \log F(\mathbf{\Lambda}_\mathbf{\Gamma}). \quad (7)$$

Note once again that  $\mathbf{\Lambda}_\mathbf{\Gamma} \equiv \mathbf{\Lambda}_\mathbf{\Gamma}(\mathbf{X}_i)$  and  $\mathbf{V}_\mu \equiv \mathbf{V}(\mu(\mathbf{X}_i))$ . If  $\mathbf{\Lambda}_\mathbf{\Gamma}$  were to be fixed as in [65], the term on the right would have no effect and minimizing that loss would correspond to optimizing the Bingham log-likelihood. To ensure  $0 \geq \lambda_1 \geq \dots \geq \lambda_{d-1}$ , we parameterize  $\boldsymbol{\lambda}$  by  $\lambda_1$  and the positive offsets  $e_2, \dots, e_{d-1}$  such that  $\lambda_k = \lambda_{k-1} - e_k$  where  $k = 2, \dots, d-1$ . This allows us to make an ordered prediction from the network.

**Extension to finite Bingham Mixture Models (BMM)** Ambiguities present in the data requires us to take into account the multimodal nature of the posterior. To achieve this, we now extend the aforementioned model to Bingham Mixture Models [67]. For the finite case, we use  $K$  different components associated with  $K$  mixture weights  $\pi_j(\mathbf{X}, \mathbf{\Gamma})$  for  $j = 1, \dots, K$ . With each component being a Bingham distribution, we can describe the density function as

$$P_\mathbf{\Gamma}(\mathbf{q} | \mathbf{X}) = \sum_{j=1}^K \pi_j(\mathbf{X}, \mathbf{\Gamma}) p_{\mathbf{\Gamma}j}(\mathbf{q} | \mathbf{X}), \quad (8)$$

where  $p_{\mathbf{\Gamma}j}(\mathbf{q} | \mathbf{X})$  are the  $K$  component distributions and  $\pi_j(\mathbf{X}, \mathbf{\Gamma})$  the mixture weights s.t.  $\sum_j \pi_j(\mathbf{X}, \mathbf{\Gamma}) = 1$ . The model can again be trained by maximizing the log-likelihood, but this time of the mixture model [79, 83]

$$\mathbf{\Gamma}^* = \arg \max_{\mathbf{\Gamma}} \log \mathcal{L}_m(\mathbf{\Gamma}|\mathcal{X}, \mathbf{Q}) \quad (9)$$

$$\log \mathcal{L}_m(\mathbf{\Gamma}|\mathcal{X}, \mathbf{Q}) = \sum_{i=1}^N \log \sum_{j=1}^K \pi_j(\mathbf{X}_i, \mathbf{\Gamma}) p_{\mathbf{\Gamma}j}(\mathbf{q}_i | \mathbf{X}_i). \quad (10)$$

## 5 Deeply modeling $\mu(\cdot)$ and $\Lambda(\cdot)$

Following up on the recent advances, we jointly model  $\mu(\cdot)$  and  $\Lambda(\cdot)$  by a deep residual network [39].  $\mathbf{\Gamma}$  denotes the entirety of the trainable parameters. On the output we have **fourteen** quantities per density: four for the mode quaternion, three for translation, three for  $\mathbf{\Lambda}$  the Bingham concentration, three for variances of the multivariate Gaussian and one for the weight  $\pi_j(\cdot)$ . In total our  $K$  mixture components result in  $K \times 14$  output entities. Our architecture is shown in Fig. 2 and we provide further details in the suppl. document. While a typical way to train our network is through simultaneously regressing the output variables, this is known to severely harm the accuracy [69]. Instead we exploit modern approaches to training in presence of ambiguities as we detail in what follows.

**MHP training scheme** Due to the increased dimensionality, in practice training our variational network in an unconstrained manner is likely to suffer from mode collapse, where all the heads concentrate around the same prediction. To avoid this and obtain a diverse set of modes, instead of training all branches equally by maximizing the log-likelihood of the mixture model, we follow the multi-hypotheses schemes of [53, 69] and train our model using a WTA loss function, for each branch maximizing the log-likelihood of a unimodal distribution,

$$\mathbf{\Gamma}^* = \arg \max_{\mathbf{\Gamma}} \sum_{i=1}^N \sum_{j=1}^K w_{ij} \log \mathcal{L}_u(\mathbf{\Gamma} | \mathbf{X}_i, \mathbf{q}_i), \quad (11)$$

according to the associated weights  $w_{ij}$  for each of the  $k$  hypotheses. In this work, we compute the weights  $w_{ij}$  during training following RWTA [69] as

$$w_{ij} = \begin{cases} 1 - \epsilon, & \text{if } j = \arg \min_k |\mathbf{q}_i - \hat{\mathbf{q}}_{ik}|, \\ \frac{\epsilon}{K-1}, & \text{otherwise} \end{cases}, \quad (12)$$

where  $\hat{\mathbf{q}}_{ik}$  is the predicted mode of a single Bingham distribution. Note that WTA [37] would amount to updating only the branch of the best hypothesis and EWTA [53] the top  $k$  branches closest to the ground truth. However, for our problem, we found RWTA to be a more reliable machinery. Finally, to obtain the desired continuous distribution, we train the weights of our Bingham mixture model using the following loss function:

$$\mathcal{L}_{\pi}(\mathbf{\Gamma} | \mathcal{X}, \mathbf{Q}) = \sum_{i=1}^N \sum_{j=1}^K \sigma(\hat{\pi}_j(\mathbf{X}_i, \mathbf{\Gamma}), y_{ij}), \quad (13)$$

where  $\sigma$  is the cross-entropy,  $\hat{\pi}(\mathbf{X}, \mathbf{\Gamma})$  the predicted weight of the neural network and  $y_{ij}$  the associated label of the mixture model component given as

$$y_{ij} = \begin{cases} 1, & \text{if } j = \arg \min_k |\mathbf{q}_i - \hat{\mathbf{q}}_{ik}|, \\ 0, & \text{otherwise} \end{cases}. \quad (14)$$

Our final loss, therefore, consists of the weighted likelihood for a unimodal distribution of each branch and the loss of our mixture weights,  $\mathcal{L}_{\pi}(\mathbf{\Gamma} | \mathcal{X}, \mathbf{Q})$ .

**Incorporating translations** We model translations  $\{\mathbf{t}_i \in \mathbb{R}^3\}_i$  by the standard Gaussian distributions with covariances  $\{\mathbf{\Sigma}_i \in \mathbb{R}^{3 \times 3} \succeq 0\}_i$ . Hence, we use the ordinary MDNs [9] to handle them. Yet, once again, during training we apply the MHP scheme explained above to avoid mode collapse and diversify the predictions. In practice, we first train the network to predict the translation and its variance. Then, intuitively, recovering the associated rotation should be an easier task, after which we fine-tune the network on all components of the distribution. Such split has already been shown to be prosperous in prior work [23].



Table 1: Evaluation in non-ambiguous scenes, displayed is the median rotation and translation error. (Numbers for MapNet on the Cambridge Landmarks dataset are taken from [74]). BPN depicts Bayesian-PoseNet [14]. *Uni* and *BMDN* refer to our unimodal version and Bingham-MDN respectively.

Dataset [° / m]	7-Scenes							Cambridge Landmarks				
	Chess	Fire	Heads	Office	Pumpkin	Kitchen	Stairs	Kings	Hospital	Shop	St. Marys	Street
PoseNet	4.48/0.13	<b>11.3</b> /0.27	13.0/0.17	5.55/0.19	4.75/0.26	5.35/0.23	12.4/0.35	<b>1.04</b> /0.88	<b>3.29</b> /3.2	<b>3.78</b> /0.88	<b>3.32</b> /1.57	25.5/20.3
MapNet	<b>3.25</b> / <b>0.08</b>	11.69/0.27	13.2/0.18	<b>5.15</b> / <b>0.17</b>	<b>4.02</b> / <b>0.22</b>	<b>4.93</b> /0.23	12.08/0.3	1.89/1.07	3.91/1.94	4.22/1.49	4.53/2.0	-
BPN	7.24/0.37	13.7/0.43	<b>12.0</b> /0.31	8.04/0.48	7.08/0.61	7.54/0.58	13.1/ 0.48	4.06/1.74	5.12/2.57	7.54/1.25	8.38/2.11	-
VidLoc	-/0.18	-/ <b>0.26</b>	-/0.14	-/0.26	-/0.36	-/0.31	-/ <b>0.26</b>	-	-	-	-	-
Uni	4.97/0.1	12.87/0.27	14.05/ <b>0.12</b>	7.52/0.2	7.11/0.23	8.25/ <b>0.19</b>	13.1/0.28	1.77/0.88	3.71/ <b>1.93</b>	4.74/ <b>0.8</b>	6.19/1.84	<b>24.1</b> /16.8
BMDN	4.35/0.1	11.86/0.28	12.76/ <b>0.12</b>	6.55/0.19	6.9/ <b>0.22</b>	8.08/0.21	<b>9.98</b> /0.31	2.08/ <b>0.83</b>	3.64/2.16	4.93/0.92	6.03/ <b>1.37</b>	36.9/ <b>9.7</b>

**Inference** Rather than reporting the conditional average which can result in label blur, we propose to obtain a single best estimate according to the weighted mode, where we choose the best mixture component according to its mixture weight and pick the mode as a final prediction.

We finally measure the uncertainty of the prediction according to the entropy of the resulting Bingham and Gaussian distributions, given as

$$H_B = \log F - \mathbf{\Lambda} \frac{\nabla F(\mathbf{\Lambda})}{F}, \quad \text{and} \quad H_G = \frac{c}{2} + \frac{c}{2} \log(2\pi) + \frac{1}{2} \log(|\mathbf{\Sigma}|), \quad (15)$$

respectively, where  $c = 3$  the dimension of the mean vector of the Gaussian. For a given image we first normalize the entropy values over all pose hypotheses, and finally obtain a measure of (un)certainty as the sum of both rotational ( $H_B$ ) and translational ( $H_G$ ) normalized entropy.

**Implementation details** We implement our method in Python using PyTorch library [61]. Following the current state-of-the-art direct camera pose regression methods, we use a *ResNet-34* [39] as our backbone network architecture, followed by fully-connected layers for rotation and translation, respectively. The predicted quaternions are normalized during training. We provide further details of training in the supplementary material.

## 6 Experimental Evaluation

When evaluating our method we consider two cases: (1) camera relocalization in non-ambiguous scenes, where our aim is to not only predict the camera pose, but the posterior of both rotation and translation that can be used to associate each pose with a measure of uncertainty; (2) we create a highly ambiguous environment, where similar looking images are captured from very different viewpoints. We show the problems current regression methods suffer from in handling such scenarios and in contrast show the merit of our proposed method.

**Error metrics** Note that, under ambiguities a best mode is unlikely to exist. In those cases, as long as we can generate a hypothesis that is close to the Ground Truth (GT), our network is considered successful. For this reason, in addition to the standard metrics and the weighted mode, we will also speak of the so called *Oracle* error, assuming an oracle that is able to choose the best of all predictions: the one closest to the GT. In addition, we report the *Self-EMD* (SEMD) [53], the earth movers distance [68] of turning a multi-modal distribution into a unimodal one. With this measure we can evaluate the diversity of predictions, where the unimodal distribution is chosen as the predicted mode of the corresponding method. Note that this measure by itself does not give any indication about the accuracy of the prediction.

**Datasets** In addition to the standard datasets of 7-Scenes [76] and Cambridge Landmarks [46], we created synthetic as well as real datasets, that are specifically designed to contain repetitive structures and allow us to assess the real benefits of our approach. For synthetic data we render table models from 3DWarehouse<sup>1</sup> and create camera trajectories, e.g. a circular movement around the object, such that ambiguous views are ensured to be included in our dataset. Specifically we use a *dining table* and a *round table* model with discrete modes of ambiguities. In addition, we create highly ambiguous real scenes using Google Tango and the graph-based SLAM approach RTAB-Map [52]. We acquire RGB and depth images as well as distinct ground truth camera trajectories for training and testing. We also reconstruct those scenes. However, note that only the RGB images and corresponding camera poses are required to train our model and the reconstructions are used for visualization only. In particular, our training and test sets consist of 2414 and 1326 frames, respectively. Note that our network sees a single pose label per image. We provide further details, visualizations and evaluations in our supplementary material.

**Baselines and SoTA** We compare our approach to current state-of-the-art direct camera pose regression methods, PoseNet [44] and MapNet [14], that output a single pose prediction. More importantly, we assess our performance against two state-of-the-art approaches, namely BayesianPoseNet [43] and VidLoc [20], that are most related to our work and predict a distribution over the pose space by using dropout and mixture density networks, respectively. We further include the *unimodal* predictions as well as BMMs trained using mixture density networks [9, 31] as baselines. We coin the latter Bingham-MDN or in short *BMDN*.

## 6.1 Evaluation in non-ambiguous scenes

We first evaluate our method on the publicly available 7-Scenes [76] and Cambridge Landmarks [46] datasets. As most of the scenes contained in these datasets do not show highly ambiguous environments, we consider them to be non-ambiguous. Though, we can not guarantee that some ambiguous views might

<sup>1</sup> <https://3dwarehouse.sketchup.com/>

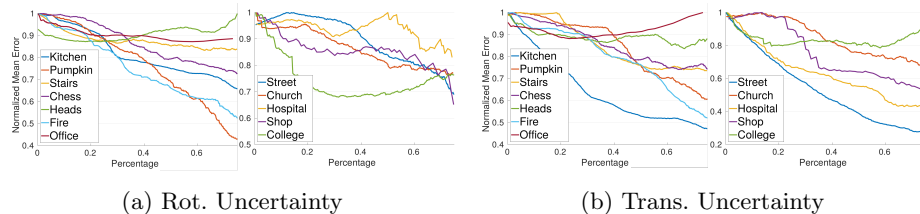


Fig. 3: Uncertainty evaluation on the 7-Scenes and Cambridge Landmarks datasets, showing the correlation between predicted uncertainty and pose error. Based on the entropy of our predicted distribution uncertain samples are gradually removed. We observe that as we remove the uncertain samples the overall error drops indicating a strong correlation between our predictions and the actual erroneous estimations.

arise in these datasets as well, such as in the *Stairs* scene of the 7-Scenes dataset. Both datasets have extensively been used to evaluate camera pose estimation methods. Following the state-of-the-art, we report the median rotation and translation errors, the results of which can be found in Tab. 1. In comparison to methods that output a single pose prediction (*e.g.* PoseNet [44] and MapNet [14]), our methods achieves similar results. Yet, our network provides an additional piece of information that is the uncertainty. On the other hand, especially in translation our method outperforms uncertainty methods, namely BayesianPoseNet [43] and VidLoc [20], on most scenes.

**Uncertainty evaluation** One benefit of our method is that we can use the resulting variance of the predicted distribution as a measure of uncertainty in our predictions. The resulting correlation between pose error and uncertainty can be seen in Fig. 3, where we gradually remove the most uncertain predictions and plot the mean error for the remaining samples. The inverse correlation between the actual errors vs our confidence shows that whenever our algorithm labels a prediction as uncertain it is also likely to be a bad estimate.

It has been shown that current direct camera pose regression methods still have difficulties in generalizing to views that differ significantly from the camera trajectories seen during training [74]. However, we chose to focus on another problem these methods have to face and analyze the performance of direct regression methods in a highly ambiguous environment. In this scenario even similar trajectories can confuse the network and easily lead to wrong predictions, for which our method proposes a solution.

## 6.2 Evaluation in ambiguous scenes

We start with quantitative evaluations on our synthetic as well as real scenes before showing qualitative results. We compare our method to PoseNet and BayesianPoseNet, which we refer to as MC-Dropout. In comparison, we replace

Table 2: Ratio of correct poses on our ambiguous scenes for several thresholds.

	Threshold	PoseNet [46]	Uni. BMDN [43]	MC-Dropout	Ours-RWTA	MC-Dropout Oracle	Ours-RWTA Oracle
Blue Chairs (A)	10° / 0.1m	0.19	0.29	0.24	<b>0.39</b>	0.35	0.40
	15° / 0.2m	0.69	0.73	0.75	0.78	<b>0.81</b>	0.90
	20° / 0.3m	<b>0.90</b>	0.86	0.80	0.88	0.82	<b>0.95</b>
Meeting Table (B)	10° / 0.1m	0.0	0.02	0.01	0.04	<b>0.05</b>	0.12
	15° / 0.2m	0.05	0.12	0.07	0.13	<b>0.28</b>	0.27
	20° / 0.3m	0.10	0.19	0.10	0.22	<b>0.39</b>	<b>0.32</b>
Staircase (C)	10° / 0.1m	0.14	0.11	0.04	0.13	<b>0.18</b>	0.19
	15° / 0.2m	0.45	0.48	0.15	0.32	<b>0.50</b>	0.53
	20° / 0.3m	0.60	0.62	0.25	0.49	<b>0.68</b>	0.70
Staircase Extended (D)	10° / 0.1m	0.07	0.06	0.06	0.02	<b>0.09</b>	0.09
	15° / 0.2m	0.31	0.26	0.21	0.14	<b>0.39</b>	0.40
	20° / 0.3m	0.49	0.41	0.32	0.31	<b>0.58</b>	<b>0.64</b>
Seminar Room (E)	10° / 0.1m	<b>0.37</b>	0.11	0.06	0.18	0.35	0.36
	15° / 0.2m	0.81	0.36	0.23	0.57	<b>0.83</b>	0.83
	20° / 0.3m	0.90	0.57	0.40	0.78	<b>0.95</b>	0.90
Average	10° / 0.1m	0.15	0.12	0.08	0.15	<b>0.20</b>	0.27
	15° / 0.2m	0.46	0.39	0.28	0.39	<b>0.56</b>	0.60
	20° / 0.3m	0.60	0.53	0.37	0.54	<b>0.68</b>	0.70

the original network architecture by a ResNet, that has been shown to improve the performance of direct camera pose regression methods [14].

**Quantitative evaluations** Due to the design of our synthetic table scenes, we know that there are two and four possible modes for each image in *dining* and *round* table scenes respectively. Hence, we analyze the predictions of our model by computing the accuracy of correctly detected modes of the true posterior. A mode is considered as found if there exists one pose hypothesis that falls into a certain rotational (5°) and translational (10% of the diameter of GT camera trajectory) threshold of it. In the dining-table, MC-Dropout obtains an accuracy of 50%, finding one mode for each image, whereas the accuracy of Ours-RWTA on average achieves 96%. On round-table, our model shows an average detection rate of 99.1%, in comparison to 24.8% of MC-Dropout.

On our real scenes, we report the recall, where a pose is considered to be correct if both the rotation and translation errors are below a pre-defined threshold. Tab. 2 shows the accuracy of our baseline methods in comparison to ours for various thresholds. Especially on our *Meeting Table* scene, it can be seen that the performance of direct camera pose regression methods that suffer from mode collapse drops significantly due to the presence of ambiguities in the scene. Thanks to the diverse mode predictions of Ours-RWTA, which is indicated by the high Oracle accuracy as well as the high SEMD shown in Tab. 3, we are able to improve upon our baseline predictions. Further, by a semi-automatic labeling procedure detailed in our suppl. material, we are able to extract GT modes for the *Blue Chairs* and *Meeting Table* scenes. This way, we can evaluate the entire set of predictions against the GT. Tab. 4 shows the percentage of correctly

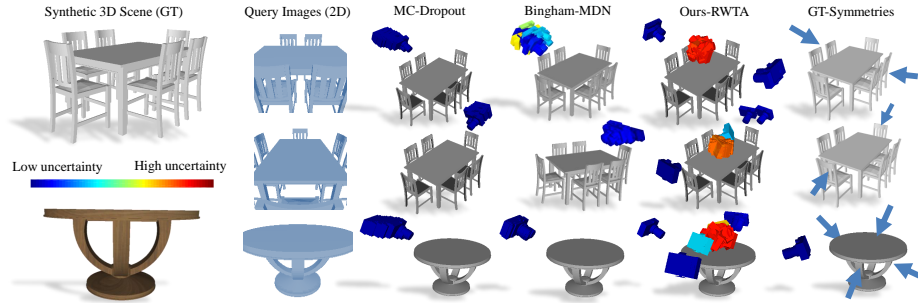


Fig. 4: Qualitative results on our synthetic *dining* and *round table* datasets. Camera poses are colored by uncertainty. Viewpoints are adjusted for best perception.

Method/Scene	A	B	C	D	E	
MC-Dropout	0.06	0.11	0.13	0.26	0.10	(a)
Ours-RWTA	<b>1.19</b>	<b>2.13</b>	<b>2.04</b>	<b>3.81</b>	<b>1.70</b>	(b)

Table 3: SEMD of our method and MC-Dropout indicating highly diverse predictions by our method in comparison to the baseline.

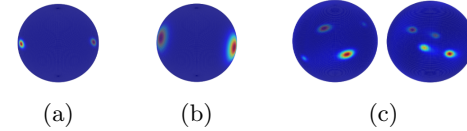


Fig. 5: Bingham distributions plotted on the unit sphere: (a) low uncertainty, (b) higher uncertainty and (c) the mixtures of Ours-RWTA.

detected modes for our method in comparison to MC-Dropout when evaluating with these GT modes. The results support our qualitative observations, that MC-Dropout suffers from mode collapse such that even with increasing threshold, the number of detected modes does not increase significantly.

**Qualitative evaluations** Qualitative results of our proposed model on our synthetic table datasets are shown in Fig. 4. MC-Dropout as well as our finite mixture model, *Bingham-MDN*, suffer from mode collapse. In comparison, the proposed MHP model is able to capture plausible, but diverse modes as well as associated uncertainties.

In contrast to other methods that obtain an uncertainty value for one prediction, we obtain uncertainty values for each hypothesis. This way, we could easily remove non-meaningful predictions, that for example can arise in the WTA and RWTA training schemes. Resulting predicted Bingham distributions are visualized in Fig. 5, by marginalizing over the angle component.

Table 4: Ratio of correctly detected modes for various translational thresholds (in meters). A and B denote *Blue Chairs* and *Meeting Table* scenes.

Scene	Method	0.1	0.2	0.3	0.4
A	MC-Dropout	0.11	0.15	0.16	0.16
	Ours-RWTA	<b>0.36</b>	<b>0.79</b>	<b>0.80</b>	<b>0.80</b>
B	MC-Dropout	0.04	0.07	0.09	0.11
	Ours-RWTA	<b>0.10</b>	<b>0.43</b>	<b>0.63</b>	<b>0.73</b>

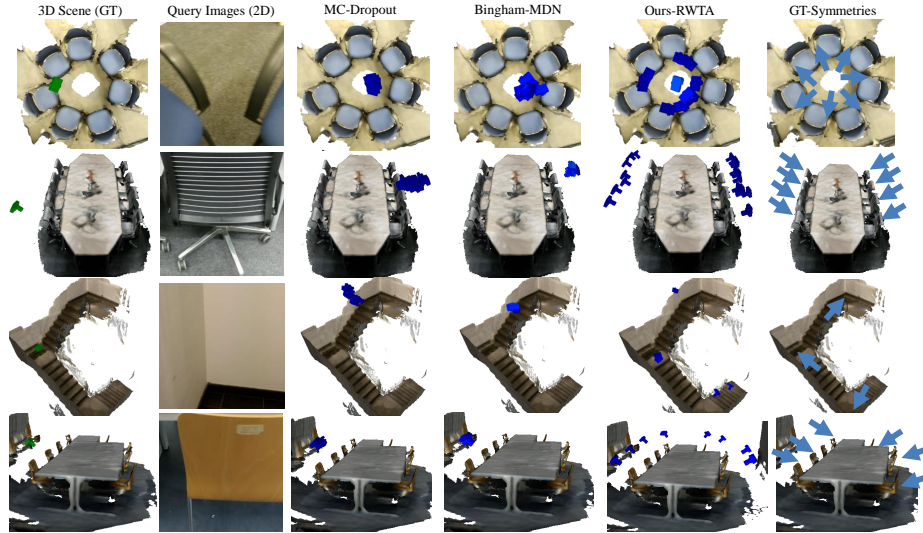


Fig. 6: Qualitative results in our ambiguous dataset. For better visualization, camera poses have been pruned by their uncertainty values.

Fig. 6 shows qualitative results on our ambiguous real scenes. Again, MC-Dropout and Bingham-MDN suffer from mode collapse. Moreover, these methods are unable to predict reasonable poses given highly ambiguous query images. This is most profound in our *Meeting Table* scene, where due to its symmetric structure the predicted camera poses fall on the opposite side of the GT one.

## 7 Conclusion

We have presented a novel method dealing with problems of direct camera pose regression in highly ambiguous environments where a unique solution to the 6DoF localization might be nonexistent. Instead, we predict camera pose hypotheses as well as associated uncertainties that finally produce a mixture model. We use the Bingham distribution to model rotations and multivariate Gaussian distribution to obtain the position of a camera. In contrast to other methods like MC-Dropout [43] or mixture density networks our training scheme is able to avoid mode collapse. Thus, we can obtain better mode predictions and improve upon the performance of camera pose regression methods in ambiguous environments while retaining the performance in non-ambiguous ones.

**Acknowledgements:** This project is supported by Bavaria California Technology Center (BaCaTeC), Stanford-Ford Alliance, NSF grant IIS-1763268, Vannevar Bush Faculty Fellowship, Samsung GRO program, the Stanford SAIL Toyota Research, and the PRIME programme of the German Academic Exchange Service (DAAD) with funds from the German Federal Ministry of Education and Research (BMBF).

## References

1. Arun Srivatsan, R., Xu, M., Zevallos, N., Choset, H.: Probabilistic pose estimation using a bingham distribution-based linear filter. *The International Journal of Robotics Research* **37**(13-14), 1610–1631 (2018) [1](#)
2. Balntas, V., Li, S., Prisacariu, V.: Relocnet: Continuous metric learning relocalisation using neural nets. In: *Proceedings of the European Conference on Computer Vision (ECCV)*. pp. 751–767 (2018)
3. Barfoot, T.D., Furgale, P.T.: Associating uncertainty with three-dimensional poses for use in estimation problems. *IEEE Transactions on Robotics* **30**(3) (2014) [5](#)
4. Bingham, C.: An antipodally symmetric distribution on the sphere. *The Annals of Statistics* pp. 1201–1225 (1974) [2](#), [4](#)
5. Birdal, T., Arbel, M., Şimşekli, U., Guibas, L.: Synchronizing probability measures on rotations via optimal transport. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition* (2020) [5](#)
6. Birdal, T., Bala, E., Eren, T., Ilic, S.: Online inspection of 3d parts via a locally overlapping camera network. In: *2016 IEEE Winter Conference on Applications of Computer Vision (WACV)*. pp. 1–10. IEEE (2016) [1](#)
7. Birdal, T., Simsekli, U.: Probabilistic permutation synchronization using the riemannian structure of the birkhoff polytope. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. pp. 11105–11116 (2019) [1](#)
8. Birdal, T., Simsekli, U., Eken, M.O., Ilic, S.: Bayesian pose graph optimization via bingham distributions and tempered geodesic mcmc. In: *Advances in Neural Information Processing Systems*. pp. 308–319 (2018) [1](#), [5](#)
9. Bishop, C.M.: Mixture density networks (1994) [3](#), [4](#), [8](#), [10](#)
10. Bourmaud, G., Mégret, R., Arnaudon, M., Giremus, A.: Continuous-discrete extended kalman filter on matrix lie groups using concentrated gaussian distributions. *Journal of Mathematical Imaging and Vision* **51**(1), 209–228 (2015) [5](#)
11. Brachmann, E., Krull, A., Nowozin, S., Shotton, J., Michel, F., Gumhold, S., Rother, C.: Dsac-differentiable ransac for camera localization. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition* (2017) [3](#)
12. Brachmann, E., Michel, F., Krull, A., Ying Yang, M., Gumhold, S., et al.: Uncertainty-driven 6d pose estimation of objects and scenes from a single rgb image. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. pp. 3364–3372 (2016) [3](#)
13. Brachmann, E., Rother, C.: Learning less is more-6d camera localization via 3d surface regression. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. pp. 4654–4662 (2018) [3](#)
14. Brahmbhatt, S., Gu, J., Kim, K., Hays, J., Kautz, J.: Geometry-aware learning of maps for camera localization. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. pp. 2616–2625 (2018) [3](#), [9](#), [10](#), [11](#), [12](#)
15. Breiman, L.: Random forests. *Machine learning* **45**(1), 5–32 (2001) [3](#)
16. Bui, M., Albarqouni, S., Ilic, S., Navab, N.: Scene coordinate and correspondence learning for image-based localization. In: *British Machine Vision Conference (BMVC)* (2018) [3](#)
17. Busam, B., Birdal, T., Navab, N.: Camera pose filtering with local regression geodesics on the riemannian manifold of dual quaternions. In: *Proceedings of the IEEE International Conference on Computer Vision Workshops*. pp. 2436–2445 (2017) [5](#)



18. Cadena, C., Carlone, L., Carrillo, H., Latif, Y., Scaramuzza, D., Neira, J., Reid, I., Leonard, J.J.: Past, present, and future of simultaneous localization and mapping: Toward the robust-perception age. *IEEE Transactions on robotics* **32**(6) (2016) [1](#)
19. Cayley, A.: Sur quelques propriétés des déterminants gauches. *Journal für die reine und angewandte Mathematik* **32**, 119–123 (1846)
20. Clark, R., Wang, S., Markham, A., Trigoni, N., Wen, H.: Vidloc: A deep spatio-temporal model for 6-dof video-clip relocalization. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition* (2017) [4](#), [10](#), [11](#)
21. Corona, E., Kundu, K., Fidler, S.: Pose estimation for objects with rotational symmetry. In: *2018 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. pp. 7215–7222. IEEE (2018) [4](#)
22. Cui, H., Radosavljevic, V., Chou, F.C., Lin, T.H., Nguyen, T., Huang, T.K., Schneider, J., Djuric, N.: Multimodal trajectory predictions for autonomous driving using deep convolutional networks. In: *2019 International Conference on Robotics and Automation (ICRA)*. pp. 2090–2096. IEEE (2019) [3](#)
23. Deng, H., Birdal, T., Ilic, S.: 3d local features for direct pairwise registration. In: *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)* (June 2019) [8](#)
24. Deng, J., Dong, W., Socher, R., Li, L.J., Li, K., Fei-Fei, L.: Imagenet: A large-scale hierarchical image database. In: *2009 IEEE conference on computer vision and pattern recognition*. pp. 248–255. Ieee (2009)
25. Durrant-Whyte, H., Bailey, T.: Simultaneous localization and mapping: part i. *IEEE robotics & automation magazine* **13**(2), 99–110 (2006) [1](#)
26. Falorsi, L., de Haan, P., Davidson, T.R., Forré, P.: Reparameterizing distributions on lie groups. *arXiv preprint arXiv:1903.02958* (2019) [5](#)
27. Feng, W., Tian, F.P., Zhang, Q., Sun, J.: 6d dynamic camera relocalization from single reference image. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. pp. 4049–4057 (2016) [3](#)
28. Firman, M., Campbell, N.D., Agapito, L., Brostow, G.J.: Diversenet: When one right answer is not enough. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. pp. 5598–5607 (2018) [3](#)
29. Fischler, M.A., Bolles, R.C.: Random sample consensus: a paradigm for model fitting with applications to image analysis and automated cartography. *Communications of the ACM* **24**(6), 381–395 (1981) [3](#)
30. Gal, Y., Ghahramani, Z.: Dropout as a bayesian approximation: Representing model uncertainty in deep learning. In: *international conference on machine learning*. pp. 1050–1059 (2016) [3](#), [4](#)
31. Gilitschenski, I., Sahoo, R., Schwarting, W., Amini, A., Karaman, S., Rus, D.: Deep orientation uncertainty learning based on a bingham loss. In: *International Conference on Learning Representations* (2020) [10](#)
32. Glover, J., Kaelbling, L.P.: Tracking the spin on a ping pong ball with the quaternion bingham filter. In: *2014 IEEE International Conference on Robotics and Automation (ICRA)*. pp. 4133–4140 (May 2014) [5](#)
33. Glover, J., Bradski, G., Rusu, R.B.: Monte carlo pose estimation with quaternion kernels and the bingham distribution. In: *Robotics: science and systems* (2012) [5](#)
34. Glover, J.M.: The quaternion Bingham distribution, 3D object detection, and dynamic manipulation. Ph.D. thesis, Massachusetts Institute of Technology (2014) [5](#)
35. Grassia, F.S.: Practical parameterization of rotations using the exponential map. *Journal of graphics tools* **3**(3), 29–48 (1998) [5](#)



36. Guo, C., Pleiss, G., Sun, Y., Weinberger, K.Q.: On calibration of modern neural networks. In: Proceedings of the 34th International Conference on Machine Learning-Volume 70. pp. 1321–1330. JMLR. org (2017) [3](#)
37. Guzman-Rivera, A., Batra, D., Kohli, P.: Multiple choice learning: Learning to produce multiple structured outputs. In: Advances in Neural Information Processing Systems. pp. 1799–1807 (2012) [3](#), [8](#)
38. Haarbach, A., Birdal, T., Ilic, S.: Survey of higher order rigid body motion interpolation methods for keyframe animation and continuous-time trajectory estimation. In: 3D Vision (3DV), 2018 Sixth International Conference on. pp. 381–389. IEEE (2018). <https://doi.org/10.1109/3DV.2018.00051> [5](#)
39. He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recognition. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 770–778 (2016) [3](#), [7](#), [9](#)
40. Herz, C.S.: Bessel functions of matrix argument. *Annals of Mathematics* **61**(3), 474–523 (1955), <http://www.jstor.org/stable/1969810> [5](#)
41. Hinterstoisser, S., Lepetit, V., Ilic, S., Holzer, S., Bradski, G., Konolige, K., Navab, N.: Model based training, detection and pose estimation of texture-less 3d objects in heavily cluttered scenes. In: Asian conference on computer vision. Springer (2012) [1](#)
42. Horaud, R., Conio, B., Le Boulleux, O., Lacolle, B.: An analytic solution for the perspective 4-point problem. In: Proceedings CVPR’89: IEEE Computer Society Conference on Computer Vision and Pattern Recognition. IEEE (1989) [1](#)
43. Kendall, A., Cipolla, R.: Modelling uncertainty in deep learning for camera relocalization. In: 2016 IEEE international conference on Robotics and Automation (ICRA). pp. 4762–4769. IEEE (2016) [1](#), [4](#), [10](#), [11](#), [12](#), [14](#)
44. Kendall, A., Cipolla, R., et al.: Geometric loss functions for camera pose regression with deep learning. In: Proc. CVPR. vol. 3, p. 8 (2017) [10](#), [11](#)
45. Kendall, A., Gal, Y.: What uncertainties do we need in bayesian deep learning for computer vision? In: Advances in neural information processing systems (2017) [4](#)
46. Kendall, A., Grimes, M., Cipolla, R.: PoseNet: A convolutional network for real-time 6-dof camera relocalization. In: Proceedings of the IEEE international conference on computer vision. pp. 2938–2946 (2015) [3](#), [4](#), [10](#), [12](#)
47. Kingma, D.P., Ba, J.: Adam: A method for stochastic optimization. arXiv preprint arXiv:1412.6980 (2014)
48. Kingma, D.P., Welling, M.: Auto-encoding variational bayes. arXiv preprint arXiv:1312.6114 (2013) [4](#)
49. Kume, A., Wood, A.T.: Saddlepoint approximations for the bingham and fisher-bingham normalising constants. *Biometrika* **92**(2), 465–476 (2005) [5](#)
50. Kurz, G., Gilitschenski, I., Julier, S., Hanebeck, U.D.: Recursive estimation of orientation based on the bingham distribution. In: Information Fusion (FUSION), 2013 16th International Conference on. pp. 1487–1494. IEEE (2013) [5](#)
51. Kurz, G., Gilitschenski, I., Pfaff, F., Drude, L., Hanebeck, U.D., Haeb-Umbach, R., Siegwart, R.Y.: Directional statistics and filtering using libdirectional. arXiv preprint arXiv:1712.09718 (2017) [5](#)
52. Labbé, M., Michaud, F.: Rtab-map as an open-source lidar and visual simultaneous localization and mapping library for large-scale and long-term online operation. *Journal of Field Robotics* **36**(2), 416–446 (2019) [10](#)
53. Makansi, O., Ilg, E., Cicek, O., Brox, T.: Overcoming limitations of mixture density networks: A sampling and fitting framework for multimodal future prediction. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 7144–7153 (2019) [2](#), [3](#), [8](#), [10](#)

54. Manhardt, F., Arroyo, D.M., Rupperecht, C., Busam, B., Birdal, T., Navab, N., Tombari, F.: Explaining the ambiguity of object detection and 6d pose from visual data. In: International Conference of Computer Vision. IEEE/CVF (2019) [1](#), [3](#)
55. Marder-Eppstein, E.: Project tango. pp. 25–25 (07 2016)
56. Mardia, K.V., Jupp, P.E.: Directional statistics. John Wiley & Sons (2009) [4](#)
57. Massiceti, D., Krull, A., Brachmann, E., Rother, C., Torr, P.H.: Random forests versus neural networks—what’s best for camera localization? In: 2017 IEEE International Conference on Robotics and Automation (ICRA). IEEE (2017) [3](#)
58. McLachlan, G.J., Basford, K.E.: Mixture models: Inference and applications to clustering, vol. 84. M. Dekker New York (1988)
59. Morawiec, A., Field, D.: Rodrigues parameterization for orientation and misorientation distributions. Philosophical Magazine A **73**(4), 1113–1130 (1996) [5](#)
60. Murray, R.M.: A mathematical introduction to robotic manipulation. CRC press (1994) [5](#)
61. Paszke, A., Gross, S., Chintala, S., Chanan, G., Yang, E., DeVito, Z., Lin, Z., Desmaison, A., Antiga, L., Lerer, A.: Automatic differentiation in PyTorch. In: NIPS Autodiff Workshop (2017) [9](#)
62. Peretroukhin, V., Wagstaff, B., Giamou, M., Kelly, J.: Probabilistic regression of rotations using quaternion averaging and a deep multi-headed network. arXiv preprint arXiv:1904.03182 (2019) [4](#)
63. Piasco, N., Sidibé, D., Demonceaux, C., Gouet-Brunet, V.: A survey on visual-based localization: On the benefit of heterogeneous data. Pattern Recognition **74**, 90–109 (2018) [1](#)
64. Pitteri, G., Ramamonjisoa, M., Ilic, S., Lepetit, V.: On object symmetries and 6d pose estimation from images. In: 3D Vision (3DV). IEEE (2019) [4](#)
65. Prokudin, S., Gehler, P., Nowozin, S.: Deep directional statistics: Pose estimation with uncertainty quantification. In: Proceedings of the European Conference on Computer Vision (ECCV). pp. 534–551 (2018) [4](#), [6](#), [7](#)
66. Qi, C.R., Litany, O., He, K., Guibas, L.J.: Deep hough voting for 3d object detection in point clouds. In: The IEEE International Conference on Computer Vision (ICCV) (October 2019) [1](#)
67. Riedel, S., Marton, Z.C., Kriegel, S.: Multi-view orientation estimation using bingham mixture models. In: 2016 IEEE international conference on automation, quality and testing, robotics (AQTR). pp. 1–6. IEEE (2016) [7](#)
68. Rubner, Y., Tomasi, C., Guibas, L.J.: The earth mover’s distance as a metric for image retrieval. International journal of computer vision **40**(2), 99–121 (2000) [10](#)
69. Rupperecht, C., Laina, I., DiPietro, R., Baust, M., Tombari, F., Navab, N., Hager, G.D.: Learning in an uncertain world: Representing ambiguity through multiple hypotheses. In: Proceedings of the IEEE International Conference on Computer Vision. pp. 3591–3600 (2017) [2](#), [3](#), [7](#), [8](#)
70. Russakovsky, O., Deng, J., Su, H., Krause, J., Satheesh, S., Ma, S., Huang, Z., Karpathy, A., Khosla, A., Bernstein, M., Berg, A.C., Fei-Fei, L.: ImageNet Large Scale Visual Recognition Challenge. International Journal of Computer Vision (IJCV) **115**(3), 211–252 (2015). <https://doi.org/10.1007/s11263-015-0816-y>
71. Salas-Moreno, R.F., Newcombe, R.A., Strasdat, H., Kelly, P.H., Davison, A.J.: Slam++: Simultaneous localisation and mapping at the level of objects. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 1352–1359 (2013) [1](#)
72. Salimans, T., Goodfellow, I., Zaremba, W., Cheung, V., Radford, A., Chen, X.: Improved techniques for training gans. In: Advances in neural information processing systems. pp. 2234–2242 (2016) [2](#)

73. Sattler, T., Havlena, M., Radenovic, F., Schindler, K., Pollefeys, M.: Hyperpoints and fine vocabularies for large-scale location recognition. In: Proceedings of the IEEE International Conference on Computer Vision. pp. 2102–2110 (2015) [1](#)
74. Sattler, T., Zhou, Q., Pollefeys, M., Leal-Taixe, L.: Understanding the limitations of cnn-based absolute camera pose regression. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 3302–3312 (2019) [9](#), [11](#)
75. Schonberger, J.L., Frahm, J.M.: Structure-from-motion revisited. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (2016) [1](#)
76. Shotton, J., Glocker, B., Zach, C., Izadi, S., Criminisi, A., Fitzgibbon, A.: Scene coordinate regression forests for camera relocalization in rgb-d images. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 2930–2937 (2013) [1](#), [3](#), [10](#)
77. Simonyan, K., Zisserman, A.: Very deep convolutional networks for large-scale image recognition. CoRR **abs/1409.1556** (2014) [3](#)
78. Subbarao, R., Meer, P.: Nonlinear mean shift over riemannian manifolds. International journal of computer vision **84**(1), 1 (2009)
79. Suvrit, S., Ley, C., Verdebout, T.: Directional statistics in machine learning: A brief review. In: Applied Directional Statistics. Chapman and Hall/CRC (2018) [7](#)
80. Szegedy, C., Vanhoucke, V., Ioffe, S., Shlens, J., Wojna, Z.: Rethinking the inception architecture for computer vision. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 2818–2826 (2016)
81. Ullman, S.: The interpretation of structure from motion. Proceedings of the Royal Society of London. Series B. Biological Sciences **203**(1153), 405–426 (1979) [1](#)
82. Valentin, J., Niefner, M., Shotton, J., Fitzgibbon, A., Izadi, S., Torr, P.H.: Exploiting uncertainty in regression forests for accurate camera relocalization. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 4400–4408 (2015) [3](#)
83. Yamaji, A.: Genetic algorithm for fitting a mixed bingham distribution to 3d orientations: a tool for the statistical and paleostress analyses of fracture orientations. Island Arc **25**(1), 72–83 (2016) [7](#)
84. Zakharov, S., Shugurov, I., Ilic, S.: Dpod: 6d pose object detector and refiner. In: The IEEE International Conference on Computer Vision (ICCV) (2019) [1](#)
85. Zeisl, B., Sattler, T., Pollefeys, M.: Camera pose voting for large-scale image-based localization. In: Proceedings of the IEEE International Conference on Computer Vision. pp. 2704–2712 (2015) [1](#)
86. Zhou, Q.Y., Park, J., Koltun, V.: Open3D: A modern library for 3D data processing. arXiv:1801.09847 (2018)
87. Zhou, Y., Barnes, C., Lu, J., Yang, J., Li, H.: On the continuity of rotation representations in neural networks. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 5745–5753 (2019)
88. Zolfaghari, M., Çiçek, Ö., Ali, S.M., Mahdisoltani, F., Zhang, C., Brox, T.: Learning representations for predicting future activities. arXiv:1905.03578 (2019) [3](#)