

Supplementary Material:

Learning Analytical Posterior Probability for Human Mesh Recovery

In this supplementary material, we provide more explanations about distribution parameters (Sec. 1), and detailed derivations of crucial theorems (Sec. 2). We also describe the sampling procedure (Sec. 3). To demonstrate the probability models intuitively, we visualize them contrastively (Sec. 4). For some unmentioned details, we simply list them for completeness (Sec. 5).

1. Parameter explanation

1.1. matrix Fisher distribution

In this part, we try to illustrate the parameter \mathbf{F} of matrix Fisher distribution from a simple geometric perspective. As stated in the main manuscript, the probability density function of matrix Fisher distribution $\mathcal{MF}(\cdot)$ over $\mathcal{SO}(n)$ [1, 5] is as follows:

$$p(\mathbf{R}; \mathbf{F}) = \frac{1}{c(\mathbf{F})} \exp(\text{tr}(\mathbf{F}^T \mathbf{R})) \sim \mathcal{MF}(\mathbf{F}).$$

Suppose there are two sets of unit vectors $\{\mathbf{l}_i\}$ and $\{\mathbf{d}_i\}$, and the vectors in $\{\mathbf{l}_i\}$ are linearly independent. To calculate a rotation \mathbf{R} that can transform $\{\mathbf{l}_i\}$ to $\{\mathbf{d}_i\}$, the estimation can be derived by minimizing the L_2 distance:

$$\begin{aligned} \hat{\mathbf{R}} &= \arg \min_{\mathbf{R}} \sum_i \|\mathbf{d}_i - \mathbf{R}\mathbf{l}_i\|_2^2 \\ &= \arg \min_{\mathbf{R}} \sum_i (\mathbf{d}_i^T \mathbf{d}_i + \mathbf{l}_i^T \mathbf{l}_i - 2\mathbf{d}_i^T \mathbf{R}\mathbf{l}_i) \\ &= \arg \max_{\mathbf{R}} \sum_i \mathbf{d}_i^T \mathbf{R}\mathbf{l}_i = \arg \max_{\mathbf{R}} \sum_i \text{tr}(\mathbf{l}_i \mathbf{d}_i^T \mathbf{R}) \\ &= \arg \max_{\mathbf{R}} \text{tr}[\sum_i \mathbf{l}_i \mathbf{d}_i^T \mathbf{R}]. \end{aligned}$$

If we define $\mathbf{F} \triangleq s \sum_i \mathbf{d}_i \mathbf{l}_i^T$ with the scale factor s , then the above equation has a similar formulation as the standard matrix Fisher distribution. In this sense, we have:

$$\mathbf{F}' = \mathbf{F} + \kappa \mathbf{d} \mathbf{l}^T = s \sum_i (\mathbf{d}_i \mathbf{l}_i^T) + \kappa \mathbf{d} \mathbf{l}^T,$$

which means the original parameter \mathbf{F} can be comprehend as the sum of the outer product of some paired vectors $(\mathbf{l}_i, \mathbf{d}_i)$, while the updated parameter \mathbf{F}' is equivalent to fusing more paired vectors on the basis of the original data.

1.2. von Mises-Fisher distribution

In this part, we illustrate the relation between normal distribution and von Mises-Fisher distribution. As stated in the main manuscript, the probability density function of von Mises-Fisher distribution $\mathcal{VMF}(\cdot)$ [6] is as follows:

$$p(\mathbf{d}; \kappa, \mathbf{m}) = \frac{1}{c(\kappa)} \exp(\kappa \mathbf{m}^T \mathbf{d}) \sim \mathcal{VMF}(\mathbf{m}, \kappa).$$

Remark. The von Mises-Fisher distribution becomes the uniform distribution on the sphere for $\kappa = 0$, and it approximates the wrapped normal distribution with the same mean \mathbf{m} and variance κ^{-1} for a large κ :

$$\begin{aligned} \mathcal{VMF}(\mathbf{m}, 0) &= \mathcal{U}(\mathcal{S}^{n-1}), \\ \mathcal{VMF}(\mathbf{m}, \kappa) &\approx \mathcal{WN}(\mathbf{m}, \kappa^{-1}), \quad \kappa \gg 0. \end{aligned}$$

Proof. It's easy to prove the statement for $\kappa = 0$. Therefore only the case of a large κ is discussed here. Note that both \mathbf{m} and \mathbf{d} can be assumed as unit vectors, leaving their scale factors to the normalizing constant. For convenience, we represent \mathbf{m} and \mathbf{d} as trigonometric functions, i.e., $(\cos \theta_m, \sin \theta_m)$ and $(\cos \theta_d, \sin \theta_d)$, respectively. Therefore, the probability density function is as follows:

$$p(\theta_d; \kappa, \theta_m) = \frac{1}{c(\kappa)} \exp(\kappa \cdot \cos(\theta_d - \theta_m)).$$

Let $\xi = \kappa^{1/2}(\theta_d - \theta_m)$, then we can derive its probability formulation as follows:

$$\begin{aligned} p(\xi) &\propto \exp(-\kappa) \cdot \exp(\kappa \cdot \cos(\kappa^{-1/2} \xi)) \\ &\approx \exp(-\kappa + \kappa(1 - \frac{1}{2} \kappa^{-1} \xi^2)) \\ &= \exp(-\frac{1}{2} \xi^2) \sim \mathcal{N}(0, 1), \end{aligned}$$

where $\exp(-\kappa)$ is a constant term independent of θ_d for transformation, and the first two terms of Taylor series of $\cos(\cdot)$ is used. Based on the fact that $p(\xi)$ follows the standard normal distribution, we can derive that $\mathbf{d} \sim \mathcal{VMF}(\mathbf{m}, \kappa) \approx \mathcal{WN}(\mathbf{m}, \kappa^{-1})$ for large κ (the condition for Taylor expansion). Note that the wrapped normal distribution is used here due to periodicity. \square

2. Theorem derivation

In this section, we derive some theorems about the posterior probability and the corresponding property thoroughly.

2.1. Posterior distribution

Remark. For $\mathbf{l}, \mathbf{d} \in \mathbb{R}^n$, $\mathbf{R} \in \mathbb{R}^{n \times n}$, the following equation holds:

$$\text{tr}(\mathbf{l}\mathbf{d}^T \mathbf{R}) = \mathbf{l}^T \mathbf{R}^T \mathbf{d}.$$

Proof.

$$\begin{aligned} \text{left} &= \text{tr}(\mathbf{l}\mathbf{d}^T \mathbf{R}) \\ &= \sum_i (\mathbf{l}\mathbf{d}^T \mathbf{R})_{ii} = \sum_i \mathbf{l}_i (\mathbf{d}^T \mathbf{R})_i = \mathbf{l}^T (\mathbf{d}^T \mathbf{R})^T \\ &= \mathbf{l}^T \mathbf{R}^T \mathbf{d} = \text{right}. \end{aligned}$$

□

Theorem 2.1. The analytical probability of rotation $\mathbf{R} \in SO(n)$ conditioned on bone direction $\mathbf{d} \in S^{n-1}$ satisfies:

$$p(\mathbf{R}|\mathbf{d}) \sim \mathcal{MF}(\mathbf{F} + \kappa \mathbf{d}\mathbf{l}^T).$$

Proof.

$$\begin{aligned} p(\mathbf{R}|\mathbf{d}) &= \frac{p(\mathbf{R}) \cdot p(\mathbf{d}|\mathbf{R})}{p(\mathbf{d})} \propto p(\mathbf{R}) \cdot p(\mathbf{d}|\mathbf{R}) \\ &= \frac{1}{c(\mathbf{F})c(\kappa)} \cdot \exp(\text{tr}(\mathbf{F}^T \mathbf{R}) + \kappa \mathbf{l}^T \mathbf{R}^T \mathbf{d}) \\ &= \frac{1}{c} \exp(\text{tr}(\mathbf{F}^T \mathbf{R}) + \text{tr}(\kappa \mathbf{l}^T \mathbf{R}^T \mathbf{d})) \\ &= \frac{1}{c} \exp(\text{tr}(\mathbf{F}^T \mathbf{R}) + \text{tr}(\kappa \mathbf{l}\mathbf{d}^T \mathbf{R})) \\ &= \frac{1}{c} \exp(\text{tr}(\mathbf{F}^T \mathbf{R} + \kappa \mathbf{l}\mathbf{d}^T \mathbf{R})) \\ &= \frac{1}{c} \exp(\text{tr}[(\mathbf{F} + \kappa \mathbf{l}\mathbf{d}^T)^T \mathbf{R}]) \\ &\sim \mathcal{MF}(\mathbf{F} + \kappa \mathbf{d}\mathbf{l}^T). \end{aligned}$$

□

Theorem 2.2. (General form) The analytical probability of rotation $\mathbf{R} \in SO(n)$ conditioned on directional observations $\mathbf{d} \in S^{n-1}$ and rotational observations $\mathbf{D} \in SO(n)$ satisfies:

$$p(\mathbf{R}|\{\mathbf{d}_i, \mathbf{D}_j\}) \sim \mathcal{MF}(\mathbf{F} + \sum_i \kappa_i \mathbf{d}_i \mathbf{l}_i^T + \sum_j \mathbf{D}_j \mathbf{K}_j^T).$$

Proof. Similar to the bone direction, the rotational observations from other sensors can also become the observation

variables \mathbf{D} of the latent variable \mathbf{R} :

$$\begin{aligned} p(\mathbf{D}|\mathbf{R}) &= \frac{1}{c} \exp(\text{tr}[(\mathbf{R}\mathbf{K})^T \mathbf{D}]) = \frac{1}{c} \exp(\text{tr}[\mathbf{K}^T \mathbf{R}^T \mathbf{D}]) \\ &= \frac{1}{c} \exp(\text{tr}[\mathbf{D}^T \mathbf{R}\mathbf{K}]) = \frac{1}{c} \exp(\text{tr}[\mathbf{K}\mathbf{D}^T \mathbf{R}]) \\ &= \frac{1}{c} \exp(\text{tr}[(\mathbf{D}\mathbf{K}^T)^T \mathbf{R}]) \sim \mathcal{MF}(\mathbf{D}\mathbf{K}^T), \end{aligned}$$

where the parameter is decomposed into a mean term \mathbf{R} and a concentration term \mathbf{K} at first, and we use some properties of square matrices: (i) $\text{tr}(A) = \text{tr}(A^T)$; (ii) $\text{tr}(AB) = \text{tr}(BA)$. Therefore, the theorem can be derived as follows:

$$\begin{aligned} p(\mathbf{R}|\{\mathbf{d}_i, \mathbf{D}_j\}) &= \frac{p(\mathbf{R}) \cdot p(\{\mathbf{d}_i, \mathbf{D}_j\}|\mathbf{R})}{p(\{\mathbf{d}_i, \mathbf{D}_j\})} \\ &\propto p(\mathbf{R}) \cdot \prod_i p(\mathbf{d}_i|\mathbf{R}) \cdot \prod_j p(\mathbf{D}_j|\mathbf{R}) \\ &= \frac{1}{c} \exp(\text{tr}[(\mathbf{F} + \sum_i \kappa_i \mathbf{d}_i \mathbf{l}_i^T + \sum_j \mathbf{D}_j \mathbf{K}_j^T)^T \mathbf{R}]) \\ &\sim \mathcal{MF}(\mathbf{F} + \sum_i \kappa_i \mathbf{d}_i \mathbf{l}_i^T + \sum_j \mathbf{D}_j \mathbf{K}_j^T) \end{aligned}$$

□

Theorem 2.1 and 2.2 provide the proofs of Eq. (6), (9) in the main manuscript, respectively.

2.2. Property

Here we provide more details about the proof of the property (Eq. (7) in the main manuscript). First, we state a lemma for the later proof.

Lemma 2.3. (Interlacing theorem) Suppose $n \geq 2$, $\mathbf{K} \in \mathbb{R}^{n \times n}$ is Hermitian, and $\mathbf{l} \in \mathbb{R}^n$ is nonzero. Then the eigenvalues satisfy the following inequality:

$$\begin{aligned} \lambda_1(\mathbf{K}) &\leq \lambda_1(\mathbf{K} + \mathbf{l}\mathbf{l}^*) \leq \lambda_2(\mathbf{K}) \leq \dots \leq \\ \lambda_n(\mathbf{K}) &\leq \lambda_n(\mathbf{K} + \mathbf{l}\mathbf{l}^*), \end{aligned}$$

which means the eigenvalues for a Hermitian perturbation with rank 1 of a Hermitian matrix are larger than the corresponding original eigenvalues in an interlaced manner [3].

Back to the proof of the property. Since the mean term \mathbf{M} is orthogonal and satisfies $\mathbf{M}\mathbf{l} = \mathbf{d}$, we can decompose the posterior parameter \mathbf{F}' as follows:

$$\begin{aligned} \mathbf{F}' &= \mathbf{F} + \kappa \mathbf{d}\mathbf{l}^T = \mathbf{M}\mathbf{K} + \kappa \mathbf{M}\mathbf{l}\mathbf{l}^T = \mathbf{M}(\mathbf{K} + \kappa \mathbf{l}\mathbf{l}^T) \\ &= \mathbf{M}\mathbf{K}', \end{aligned}$$

where $\mathbf{K} = \mathbf{V}\Delta\mathbf{S}\mathbf{V}^T$ is a real symmetric matrix considering that both Δ and \mathbf{S} are diagonal. Besides, $\mathbf{l}\mathbf{l}^T$, the

outer product of \mathbf{l} with itself, is a symmetric matrix with rank 1. Thus \mathbf{K}' is also real symmetric, *i.e.*, Hermitian. From Lemma 2.3, we can get the conclusion that \mathbf{K}' has larger eigenvalues than \mathbf{K} , therefore the singular values of the posterior parameter \mathbf{F}' is larger than those of the prior parameter \mathbf{F} .

Fig. 1 illustrates the convergence performance regarding the accuracy of SMPL poses for two different settings, respectively, estimating prior \mathbf{F} without a keypoint branch and posterior \mathbf{F}' with a keypoint branch. The latter clearly converges faster and better.

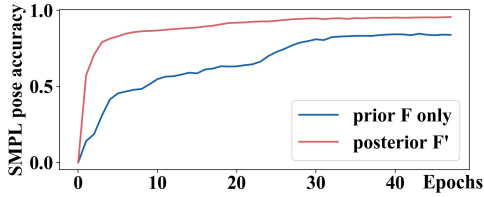


Figure 1. Convergence comparison.

3. Sampling

As we stated in the main manuscript, because of the equivalence between Bingham distribution on S^3 and matrix Fisher distribution on $\mathcal{SO}(3)$, we sample the Bingham distribution instead via differentiable rejection sampling [4, 7]. Algorithm 1 lists the detailed procedure. The proposal distribution in rejection sampling is an angular central Gaussian (ACG) distribution, the sampling of which is implemented with the reparameterization trick.

Algorithm 1: \mathcal{MF} Rejection Sampling

Input: $\mathbf{F} \in \mathbb{R}^{3 \times 3}$
Output: $\mathbf{R}_i \sim \mathcal{MF}(\mathbf{F})$

- 1 $\mathbf{U}, \mathbf{S}, \mathbf{V}^T = \text{SVD}(\mathbf{F})$, where $\mathbf{S} = \text{diag}(s_1, s_2, s_3)$
- 2 \triangleright Sample matrix Fisher via Bingham
- 3 $\mathbf{A} = \text{diag}(0, 2(s_2 + s_3), 2(s_1 + s_3), 2(s_1 + s_2))$
- 4 Solve b , s.t. $\sum_{i=1}^4 \frac{1}{b+2\mathbf{A}_i} - 1 = 0$
- 5 $\mathbf{\Omega} = \mathbf{I}_4 + \frac{2}{b}\mathbf{A}$
- 6 $M = (\frac{4}{b})^2 \exp(\frac{b-4}{2})$
- 7 \triangleright Sample Bingham via ACG as proposal distribution
- 8 **repeat**
- 9 Sample $w \sim \mathcal{U}(0, 1)$
- 10 Sample $\epsilon \sim \mathcal{N}(\mathbf{0}_4, \mathbf{I}_4)$ \triangleright reparameterization
- 11 $\mathbf{y} = (\mathbf{\Omega}^{-1})^{\frac{1}{2}} \epsilon$
- 12 Get a proposal $\mathbf{x} = \frac{\mathbf{y}}{\|\mathbf{y}\|}$ s.t. $\mathbf{x} \in S^3$
- 13 **until** $w < \frac{\exp(-\mathbf{x}^T \mathbf{A} \mathbf{x})}{M(\mathbf{x}^T \mathbf{\Omega} \mathbf{x})^{-2}}$;
- 14 $\mathbf{R}_i = \mathbf{U} \mathbf{X} \mathbf{V}^T$, where $\mathbf{X} = \text{quat_to_mat}(\mathbf{x})$
- 15 **return** \mathbf{R}_i

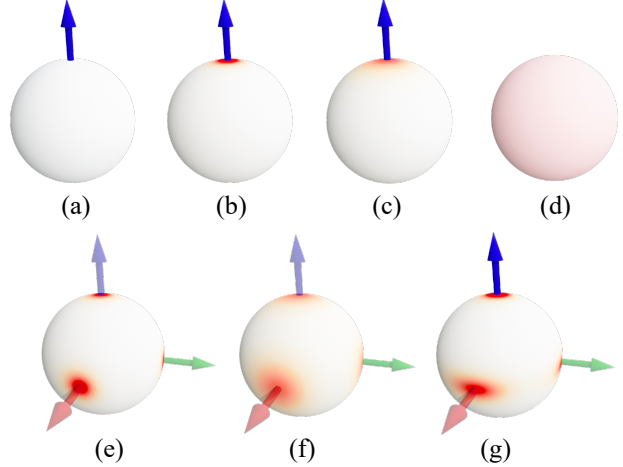


Figure 2. Comparison of different models.

4. Figure illustration

In this part, we illustrate Fig. 3 in the main manuscript thoroughly. As shown in Fig. 2, relevant spheres represent the following underlying models: (a) deterministic direction without probabilistic modeling (*e.g.*, IK estimations); (b) directional distribution with large confidence (*e.g.*, gravity sensors); (c) directional distribution with small confidence; (d) uniform distribution on the sphere; (e) rotational distribution with large confidence (*e.g.*, gyroscopes); (f) rotational distribution with small confidence; (g) rotational posterior distribution conditioned on directional information.

The eigenvalues of (e) are larger than those of (f), thus, (e) shows a more concentrated distribution. For (g), since the directional information is provided and fused, the distribution is more concentrated on the blue arrow, so the regions around the red and green arrows show oval shapes (not circles). (g) is easier to learn compared with directly regressing parameters (validated by Fig. 1), because the prior and the keypoints not only mutually narrow down the region of solutions but also have the potential to recover the ground-truth even for noisy cases.

5. Others

In this section, we provide some unmentioned details.

Error definition: Given the estimated rotation $\hat{\mathbf{R}}$, ground-truth \mathbf{R}^* , and canonical unit vector \mathbf{l} , the rotation error Err_R and direction error Err_d in the simulation experiment are as follows:

$$Err_R = \arccos\left(\frac{\text{tr}(\mathbf{R}^* \hat{\mathbf{R}}^T) - 1}{2}\right),$$

$$Err_d = \arccos(\langle \mathbf{R}^* \mathbf{l}, \hat{\mathbf{R}} \mathbf{l} \rangle).$$

Module	Block	Input	Output
Conv(R)	ResNet	Image (256, 256, 3)	$(f_c, 8, 8)$
Deconv(R)	{Deconv+BN+ReLU} $\times 3$	$(f_c, 8, 8)$	$(f_c, 64, 64)$
	1 \times 1 Conv	$(f_c, 64, 64)$	$(29 \times 64, 64, 64)$
	Soft-argmax	$(29 \times 64, 64, 64), \mathbf{s}$	$\mathbf{J}(29, 3)$
Conv(H)	HRNet	Image (256, 256, 3)	$(f_c, 64, 64)$
Deconv(H)	1 \times 1 Conv	$(f_c, 64, 64)$	$(29 \times 64, 64, 64)$
	Soft-argmax	$(29 \times 64, 64, 64), \mathbf{s}$	$\mathbf{J}(29, 3)$
MLP feature	AvgPool + Fc1 + Dropout	$(f_c, 64, 64)$	(f'_c)
	Fc2 + Dropout	(f'_c)	(f'_c)
MLP(β)	Fc(β)	(f'_c)	$\beta(10)$
MLP(\mathbf{F})	Fc(\mathbf{F})	(f'_c)	$\mathbf{F}(216)$
MLP(\mathbf{s})	Fc(\mathbf{s})	(f'_c)	$\mathbf{s}(1)$
Output	shape $\beta(10)$, parameter $\mathbf{F}(216)$, 3D keypoints $\mathbf{J}(29, 3)$		

Table 1. **Network architecture.**

Differential entropy: For the matrix Fisher distribution, the differential entropy H_F can be a representation of its uncertainty [2, 8] and calculated via Bingham distribution:

$$H_F = \log c_B - \sum_{i=1}^3 k_i \frac{\partial c_B}{\partial k_i} - \log(2\pi^2),$$

where c_B , $\mathbf{K} = [k_1, k_2, k_3]$ are the normalizing constant and the concentration parameters of Bingham distribution, respectively. We normalize the entropy from multiple sensors, since the differential entropy has not a certain range. Note that other metrics also deserve exploration.

Network architecture: Table 1 lists the structure of our network and the corresponding feature dimensions. f_c is the number of feature channel relevant to different backbones. f'_c is the MLP feature dimension. Fc1 and Fc2 are shared by MLP branches.

Model size and running time. The total numbers of trainable parameters and the running time on an RTX 2080ti GPU (batch size 1) are 27.8M / 19ms (R-34 backbone) and 73.4M / 46ms (H-48 backbone), respectively. The model size and speed are of the same level as other frameworks.

References

- [1] Thomas D Downs. Orientation statistics. *Biometrika*, 59(3):665–676, 1972. 1
- [2] Jared Marshall Glover. *The quaternion Bingham distribution, 3D object detection, and dynamic manipulation*. PhD thesis, MIT, 2014. 4
- [3] Roger A Horn and Charles R Johnson. *Matrix analysis*. Cambridge university press, 2012. 2
- [4] John T Kent, Asaad M Ganeiber, and Kanti V Mardia. A new method to simulate the bingham and related distributions in directional data analysis with applications. *arXiv*, 2013. 3
- [5] CG Khatri and Kanti V Mardia. The von mises–fisher matrix distribution in orientation statistics. *Journal of the Royal Statistical Society: Series B (Methodological)*, 39(1):95–106, 1977. 1
- [6] Kanti V Mardia, Peter E Jupp, and KV Mardia. *Directional statistics*, volume 2. Wiley Online Library, 2000. 1
- [7] Akash Sengupta, Ignas Budvytis, and Roberto Cipolla. Hierarchical kinematic probability distributions for 3d human shape and pose estimation from images in the wild. In *ICCV*, pages 11219–11229, 2021. 3
- [8] Yingda Yin, Yingcheng Cai, He Wang, and Baoquan Chen. Fishermatch: Semi-supervised rotation regression via entropy-based filtering. In *CVPR*, pages 11164–11173, 2022. 4