

Supplemental Material: Dynamic Hyperbolic Attention Network for Fine Hand-object Reconstruction

Zhiying Leng^{1,2}, Shun-Cheng Wu², Mahdi Saleh², Antonio Montanaro³, Hao Yu², Yin Wang¹, Nassir Navab², Xiaohui Liang^{1,4*}, Federico Tombari²

¹ State Key Laboratory of Virtual Reality Technology and Systems, Beihang University, China

² Computer Aided Medical Procedures, Technical University of Munich, Germany

³ Politecnico di Torino, Italy

⁴ Zhongguancun Laboratory, Beijing, China

{zhiyingleng, liang_xiaohui}@buaa.edu.cn, {shuncheng.wu, m.saleh}@tum.de, tombari@in.tum.de

Our main paper introduced a **Dynamic Hyperbolic Attention Network (DHANet)** for accurate hand-object reconstruction, which includes dynamic hyperbolic graph convolution (DHGC) and image-attention hyperbolic graph convolution (IHGC). In this supplemental document, we provide more detailed information about our method, including:

- Detailed architecture of our method in Section 7;
- More comparative experiments in Section 8;
- More reconstruction results in Section 9;

7. Detailed Architecture

7.1. Architecture for DHANet

Input: I , an 256*256 image		Object Branch	
Hand Branch		Architecture	
Components	Architecture	Components	Architecture
Encoder	ResNet-18(I)	Encoder	ResNet-18(I)
Output:	F_{sh}^E, F_{dh}^E, F_h	Output:	F_{so}^E, F_{do}^E, F_o
Decoder	FC+MANO-layer(F_h)	Decoder	AtlasNet(F_o)
Output:	m_h	Output:	m_o
DHGC	DHGC(m_h)	DHGC	DHGC(m_o)
Output:	$F_{h,mesh}^B$	Output:	$F_{o,mesh}^B$
DHGC	DHGC($F_{h,mesh}^B$)	DHGC	DHGC($F_{o,mesh}^B$)
Output:	$F_{h,mesh}^B$	Output:	$F_{o,mesh}^B$
IHGC	IHGC($F_{h,mesh}^B, F_{sh}^E, F_{dh}^E, F_{do}^E$)	IHGC	IHGC($F_{o,mesh}^B, F_{so}^E, F_{do}^E, F_{do}^E$)
Output:	$F_{o \rightarrow h}^E$	Output:	$F_{h \rightarrow o}^E$
Decoder	FC+MANO-layer($F_h, F_{o \rightarrow h}^E$)	Decoder	AtlasNet($F_o, F_{h \rightarrow o}^E$)
Output:	m_h	Output:	m_o

Table 5. The architecture of our DHANet. “E” refers to Euclidean space. “B” refers to hyperbolic space.

In this section, we introduce the detailed network architecture for our DHANet. Given an input image I with size 256×256 , DHANet reconstructs a hand mesh m_h of size 778×3 , and an object mesh m_o of size 642×3 . As listed in

*corresponding author

Table 5, our DHANet is a two-branch network, one for hand reconstruction and the other one for object reconstruction.

7.2. Architecture for DHGC

In this section, the detailed architecture of DHGC is stated. DHGC consists of projection, graph construction, and hyperbolic graph convolution. Given a mesh in $l-1$ -th layer, $v^{l-1,E}$, the l -th DHGC layer learns mesh features, $v^{l,B}$. The detailed implementation is listed in Table 6.

Components	Architecture
Exp-projection	geopt.manifold.expmat0($v^{l-1,E}$)
Build graphs	geopt_layers.KNN(k=20)
Transformation	poincare.MobiusLinear(3,64)
Activation	poincare.RadialNd(nn.LeakyReLU)
Aggregation	poincare.math.poincare_mean()
Output: $v^{l,B}$	

Table 6. The detailed architecture of DHGC. We implemented DHGC based on the geopt tool [2].

7.3. Architecture for IHGC

In this section, the detailed structure of IHGC is introduced. Given image features and mesh features, IHGC learns multi-modal features. Table 7 lists the detailed implementation of IHGC in the hand branch. The architecture of IHGC in the object branch is similar to the hand branch.

8. More comparative experiments

Comparison with spectral convolutions. Our DHANet is a spatial-based graph convolution defined in hyperbolic space. Nevertheless, we also itemize the comparison with a spectral-based graph convolution (Chebyshev graph convolution) in Table 8.

Components	Architecture
Transformation for shallow image features	$F_{sh}^E = \text{nn.Conv1d}(64, 32, \text{kernel_size}=1)$
Exp-projection	$F_{sh}^B = \text{geopt.manifold.expmap0}$
Build graphs	$f_2 = \text{geopt.layers.KNN}(k=20)$
Features transformation	$V = \text{poincare.MobiusLinear}(64, 32)(\text{concat}(f_1, f_2))$
Log-projection	$V = \text{geopt.manifold.logmap0}$
Transformation for deep image features of hand	$F_{dh}^E = \text{nn.Conv1d}(64, 32, \text{kernel_size}=1)$
Exp-projection	$F_{dh}^B = \text{geopt.manifold.expmap0}$
Build graphs	$Q = \text{geopt.layers.KNN}(k=20)$
Log-projection	$Q = \text{geopt.manifold.logmap0}$
Transformation for deep image features of object	$F_{do}^E = \text{nn.Conv1d}(64, 32, \text{kernel_size}=1)$
Exp-projection	$F_{do}^B = \text{geopt.manifold.expmap0}$
Build graphs	$K = \text{geopt.layers.KNN}(k=20)$
Log-projection	$K = \text{geopt.manifold.logmap0}$
Attention	$\text{att_weight} = \text{torch.matmul}(Q, K) / d$ $\text{att_weight} = \text{softmax}(\text{att_weight})$
Aggregation	$F_{o \rightarrow h}^E = \text{torch.matmul}(\text{att_weight}, V)$
Output: $F_{o \rightarrow h}^E$	$F_{o \rightarrow h}^E = \text{torch.sum}(F_{o \rightarrow h}^E)$

Table 7. The detailed architecture of IHGC in the hand branch. We implemented IHGC based on the geopt tool [2].

Method	Baseline [1]	Baseline+DGCNN	Baseline+ChebConv	ours
Hand error	28.1	25.9	27.6	23.8
Obj error	1579.2	1490.6	1489.8	1236.0

Table 8. Comparisons with different graph convolutions on FHB⁻.

Comparison on different datasets. Our method achieves the state-of-the-art result on FHB⁻ dataset, while yields the lower accuracy than [1] on Obman dataset. The reason for that is dataset difference: Obman dataset is a synthetic dataset with complex and various backgrounds, while FHB⁻ dataset is a real-world dataset with a simplistic kitchen environment background. This results in shallow image features learned on Obman containing more irrelevant background information, as shown in Figure 8. The more irrelevant features are, the more loose their projections are in hyperbolic space, and vice versa, as shown in Figure 8. Those loose shallow image features in hyperbolic space are unfavorable to searching the neighborhood between those features and mesh features. This causes our approach to work slightly worse in Obman dataset.

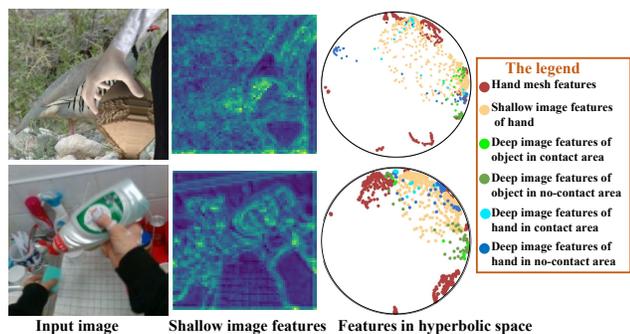


Figure 8. Feature differences between Obman dataset (the first row) and FHB⁻ dataset (the second row).

9. Additional Results

As mentioned in Section 4.3 of the main paper, image features projected to hyperbolic space preserve the spatial relationships between hand and object in images, which

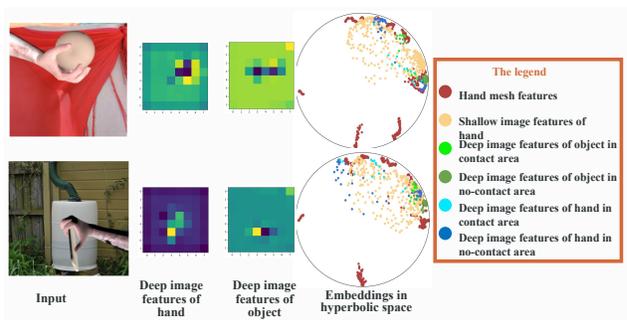


Figure 9. More visualizations for embeddings in hyperbolic space.

is beneficial for modeling the hand-object interaction. We have shown some visualizations in Fig. 7 in the main paper. This supplemental material provides more visualizations for embeddings in hyperbolic space, as shown in Fig. 9. We observed that the light blue points are close to the light green points while other points are far away. This distribution is consistent with the spatial relationship between hand and object in the image.

References

- [1] Yana Hasson, Gul Varol, Dimitrios Tzionas, Igor Kalevtykh, Michael J Black, Ivan Laptev, and Cordelia Schmid. Learning joint reconstruction of hands and manipulated objects. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 11807–11816, 2019. 2
- [2] Max Kochurov, Rasul Karimov, and Serge Kozlukov. Geopt: Riemannian optimization in pytorch, 2020. 1, 2