## *Supplementary Material* for Multi-scale Matching Networks for Semantic Correspondence

Dongyang Zhao<sup>1,2</sup>, Ziyang Song<sup>1</sup>, Zhenghao Ji<sup>1</sup>, Gangming Zhao<sup>3</sup>, Weifeng Ge \*<sup>1,2</sup>, and Yizhou Yu<sup>3</sup>

<sup>1</sup>Nebula AI Group, School of Computer Science, Fudan University
<sup>2</sup>Shanghai Key Lab of Intelligent Information Processing
<sup>3</sup>Department of Computer Science, The University of Hong Kong

## A. Parameter Table

Images are resized to  $224 \times 320$  for all datasets both in training and testing.

		ResNet-50	ResNet-101	ResNeXt-101	ResNet-101-FCN			
Conv Intra-fused	conv: output channel = 21, kernel size = $1 \times 1$ , stride = 1							
Conv Cross-fused	conv: outp	ut channel = 21, kernel s	$ize = 3 \times 3$ , stride = 1, pa	dding = 1				
Deconv	transpose conv: output	channel = 21, kernel siz	$e = 4 \times 4$ , stride = 2, pade	ding = 1, bias = False				
1.54	conv_key	conv: output channel = 10, kernel size = $1 \times 1$						
	conv_value	conv: output channel = 10, kernel size = $1 \times 1$						
LJA	conv_query	conv: outp	out channel = 10, kernel	size = $1 \times 1$				
	conv_aggregate	conv: outp	put channel = $21$ , kernel	size = $1 \times 1$				
# Params	backbone	$25.6 \times 10^6$	$44.5.0 \times 10^{6}$	$88.8 \times 10^{6}$	$54.4 \times 10^6$			
	MMNet w/o backbone	$4.8 \times 10^6$	$10.3  imes 10^6$	$10.3 \times 10^6$	$10.3  imes 10^6$			
	MMNet	$30.4 \times 10^6$	$54.8 \times 10^6$	$99.1 \times 10^6$	$64.7  imes 10^6$			
FLOPs	backbone	$11.8 \times 10^{9}$	$22.4 \times 10^{9}$	$47.0 \times 10^{9}$	$127.6 \times 10^{9}$			
	MMNet w/o backbone	$3.1 \times 10^9$	$4.6 \times 10^{9}$	$4.6 \times 10^{9}$	$12.7 \times 10^9$			
	MMNet	$14.9 \times 10^9$	$27.0 \times 10^{9}$	$51.6 \times 10^{9}$	$140.3 \times 10^9$			

Table 1. Implementation details of MMNet, and numbers of parameters and FLOPs introduced by different modules. In this table, '*Conv Intra-fused*' and '*Conv Cross-fused*' denote the convolution operation in the intra-scale and cross-scale feature enhancements separately, '*Deconv*' indicates the deconvolution operation to upscale the feature maps during the cross-scale feature enhancement, and '*LSA*' is short for the local self attention module used at the end of the intra-scale feature enhancement.

**Analysis.** MMNet introduces additional parameters only in its feature enhancement module. We follow BDCN [1] to set the parameters of scale enhancement modules. Each scale enhancement module contains four  $3 \times 3$  convolution operations with dilation 1, 4, 8, 12 respectively. The output channel numbers of these convolution operations are set to 32. The output channel numbers of other convolution and deconvolution operations are set to 21 to reduce the computational cost. We compare the additional parameters and FLOPs introduced by MMNet with four different backbones including ResNet-50, ResNet-101, ResNeXt-101 and ResNet-101-FCN. The additional parameters are no larger than 25% of that of any backbone. For the FLOPs, we add at most 26.3% computational cost when using ResNet-50 as the backbone.

## **B. Additional Results**

**Results with other backbones.** We adapt our MMNet design with other backbones: VGG-16 and DeepLab-V3, the results are listed in 2. As can be seen, our model with ResNet-101-FCN backbone performs significantly best than others.

<sup>\*</sup>Corresponding author: wfge@fudan.edu.cn

	PF-PASCAL		
Methods	0.05	0.1	0.15
MMNet <sub>ResNet-101-FCN</sub>	81.1	91.6	95.9
MMNet <sub>VGG-16</sub>	69.5	82.0	88.4
MMNet <sub>DeepLabV3-ResNet101</sub>	73.3	85.3	92.3

Table 2. Experiments on different backbones. All experiments is conducted on PF-PASCAL.



Figure 1. Key-point matching results on SPair-71k dataset [3] compared with SCOT [2] and DHPF [4]. The odd rows are the source images, and the even rows are the target images. Destination key points are denoted with crosses.



Figure 2. Key-point matching results on SPair-71k dataset [3] compared with SCOT [2] and DHPF [4]. The odd rows are the source images, and the even rows are the target images. Destination key points are denoted with crosses.



Figure 3. Warped images by thin-plate splines with the predicted key point pairs on SPair-71k dataset [3] compared with SCOT [2] and DHPF [4]



Figure 4. Warped images by thin-plate splines with the predicted key point pairs on SPair-71k dataset [3] compared with SCOT [2] and DHPF [4]

## References

- Jianzhong He, Shiliang Zhang, Ming Yang, Yanhu Shan, and Tiejun Huang. Bi-directional cascade network for perceptual edge detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3828–3837, 2019.
- [2] Yanbin Liu, Linchao Zhu, Makoto Yamada, and Yi Yang. Semantic correspondence as an optimal transport problem. In *Proceedings* of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pages 4463–4472, 2020. 3, 4, 5, 6
- [3] Juhong Min, Jongmin Lee, J. Ponce, and Minsu Cho. Spair-71k: A large-scale benchmark for semantic correspondence. ArXiv, abs/1908.10543, 2019. 3, 4, 5, 6
- [4] Juhong Min, Jongmin Lee, Jean Ponce, and Minsu Cho. Learning to compose hypercolumns for visual correspondence. In *ECCV* 2020-16th European Conference on Computer Vision. Springer, 2020. 3, 4, 5, 6