# HM: Hybrid Masking for Few-Shot Segmentation (Supplementary Material)

Seonghyeon Moon<sup>1</sup>, Samuel S. Sohn<sup>1</sup>, Honglu Zhou<sup>1</sup>, Sejong Yoon<sup>2</sup>, Vladimir Pavlovic<sup>1</sup>, Muhammad Haris Khan<sup>3</sup>, and Mubbasir Kapadia<sup>1</sup>

<sup>1</sup> Rutgers University, New Jersey, USA {sm2062,samuel.sohn,honglu.zhou,vladimir,mubbasir.kapadia}@rutgers.edu <sup>2</sup> The College of New Jersey, New Jersey, USA yoons@tcnj.edu <sup>3</sup> Mohamed Bin Zayed University of Artificial Intelligence, Abu Dhabi, UAE muhammad.haris@mbzuai.ac.ae

Our method (HM) reveals improvement in mIoU on three standard benchmarks compared to the existing state-of-the-art HSNet [5] by generating features indicative of objects' contours, edges, and complex patterns through proposed hybrid masking. In this supplementary material, we further include the following materials:

- 1. Further analysis of training efficiency.
- 2. Additional qualitative results on three bench marks.
- 3. Further analysis of failure cases.

## 1 Training efficiency

Compared to the existing HSNet [5] (baesline), out method (HSNet-HM) shows improved training efficiency not only on the COCO- $20^i$  dataset [4] but also on the other two benchmark datasets, PASCAL- $5^i$  [1] and FSS-1000 [3]. Fig. 1 and Fig. 2 display that our method tends to converge faster than HSNet [5]. Table 1 shows the number of epochs it takes to get the best performing model on the validation set. In particular, the results of COCO- $20^i$  with ResNet101 [2] show that our method (HM) has a better training efficiency.

Backbone feature	Masking methods	$\begin{vmatrix} 1\\ 20^0 \end{vmatrix}$	$PASO 20^1$	$20^2$	$\frac{5^{i}}{20^{3}}$	l-shot mEpoch	$ _{20^0}$	$\begin{array}{c} \mathrm{CO} \\ 20^1 \end{array}$	CO-2 $20^2$		-shot mEpoch	$FSS-1000^i$ 1-shot Epoch
	HSNet [5]	345	433	204	244	306.5	262	249	160	295	241.5	530
ResNet50 [2]	HSNet-HM	188	117	<b>45</b>	56	101.5	41	32	32.8	<b>23</b>	35	177
ResNet101 [2]	HSNet [5]	177	185	136	199	174.3	235	251	345	355	296.5	886
	HSNet-HM	73	95	30	<b>72</b>	67.5	52	<b>27</b>	<b>14</b>	<b>14</b>	26.8	298

Table 1: Number of epochs to reach the best model.



Fig. 1: Training profiles of HSNet[5] and HSNet-HM on  $PASCAL-5^{i}$  [1]



Fig. 2: Training profiles of HSNet[5] and HSNet-HM on FSS-1000.

#### 2 Qualitative results

HM reveal better performance in mIoU over other methods on three standard benchmarks. Fig. 3,4,5 shows the qualitative results of each benchmark. HM is able to segment the target object in the query image precisely. Also, even if only a part of the object exists in the support set or query image because of occlusion, HM can segment the target object more accurately.

# 3 Failure Cases

Although HM provides significant improvements over existing best methods, it can be further improved after analyzing the failure cases shown in Fig. 6. We discuss these failure cases below:

- When the target object in the support set or query image is very small, HM was unable to segment accurately. We believe that the reason for this is that feature maps could not contain enough information about the target object.
- A problem arose when the shape of the target object of the support set and the target image of the query image were different. HM was able to detect where the target object was, but the details failed to segment.
- If there is an object with a similar texture to the target object in the query image, there is a difficulty in segmentation. For example, the textures of cows and horses were quite similar, and because of this, the trained model could not distinguish between horses and cows.



Fig. 3: HSNet-HM Qualitative result on PASCAL-5 $^i\ [1]$ 



Fig. 4: HSNet-HM Qualitative result on FSS-1000 [3]



Fig. 5: HSNet-HM Qualitative result on COCO-20  $^i$  [4]



Fig. 6: HSNet-HM Failure cases on COCO-20^i [4]

## References

- Everingham, M., Van Gool, L., Williams, C., Winn, J., Zisserman, A.: The pascal visual object classes (voc) challenge. International Journal of Computer Vision 88, 303–338 (06 2010). https://doi.org/10.1007/s11263-009-0275-4
- He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recognition. In: 2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR). pp. 770–778. IEEE Computer Society (2016). https://doi.org/10.1109/CVPR.2016.90, https://doi.org/10.1109/CVPR.2016.90
- Li, X., Wei, T., Chen, Y.P., Tai, Y.W., Tang, C.K.: Fss-1000: A 1000-class dataset for few-shot segmentation. CVPR (2020)
- Lin, T.Y., Maire, M., Belongie, S., Bourdev, L., Girshick, R., Hays, J., Perona, P., Ramanan, D., Zitnick, C.L., Dollár, P.: Microsoft coco: Common objects in context (2015)
- Min, J., Kang, D., Cho, M.: Hypercorrelation squeeze for few-shot segmentation. In: Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV). pp. 6941–6952 (October 2021)