

Coarse Retinal Lesion Annotations Refinement via Prototypical Learning

Qinji Yu¹, Kang Dang², Ziyu Zhou¹, Yongwei Chen², and Xiaowei Ding^{1,2}(✉)

¹ Shanghai Jiao Tong University, Shanghai, China

² VoxelCloud, Inc., Los Angeles, USA

Abstract. Deep-learning-based approaches for retinal lesion segmentation often require an abundant amount of precise pixel-wise annotated data. However, coarse annotations such as circles or ellipses for outlining the lesion area can be six times more efficient than pixel-level annotation. Therefore, this paper proposes an annotation refinement network to convert a coarse annotation into a pixel-level segmentation mask. Our main novelty is the application of the prototype learning paradigm to enhance the generalization ability across different datasets or types of lesions. We also introduce a prototype weighing module to handle challenging cases where the lesion is overly small. The proposed method was trained on the publicly available IDRiD dataset and then generalized to the public DDR and our real-world private datasets. Experiments show that our approach substantially improved the initial coarse mask and outperformed the non-prototypical baseline by a large margin. Moreover, we demonstrate the usefulness of the prototype weighing module in both cross-dataset and cross-class settings.

Keywords: Prototypical Learning · Retina Lesion Segmentation · Coarse Annotation Refinement.

1 Introduction

Given the growing demand for retinal screening, automatic segmentation for retinal lesions enjoys increasing clinical relevance. By answering the issue of what lesions exist in the image and where they are located, retinal lesion segmentation algorithms assist ophthalmologists in making clinical diagnoses and assessing disease severity [18]. While recent deep-learning approaches have tremendously boosted the retinal lesion segmentation accuracy [18,9,4,19], they often require abundant expert-level-accurate, pixel-wise annotated data, which requires significant time and expense to acquire. Previous studies show that coarse annotations such as circles or ellipses for outlining the lesion area can be six times more efficient than pixel-level annotation [5]. Therefore, it is essential to study novel methodologies tailored for lower-quality coarse annotations.

Q. Yu and K. Dang contribute equally to this work.

Existing works on exploiting coarse annotations can be categorized into weakly-supervised segmentation [11,20,15,10,16,1,24], and mask refinement [21,5]. Weakly-supervised segmentation methods rely on prior assumptions such as box tightness constraint [16] and image contrast constraint [15] to utilize box-level and image-level coarse annotations. A few high-quality pixel-level retinal lesion datasets such as IDRiD [12] and DDR [8] provide precise lesion boundaries. While successful, weakly-supervised segmentation does not utilize these pre-existing lesion segmentation datasets that provide rich knowledge on lesions’ exact appearance and shape. Instead of further developing weakly-supervised segmentation methods, we propose to use such datasets by training an annotation refinement model in a data-driven manner to convert a coarse annotation into a pixel-level segmentation mask. It should be emphasized that our work is significantly different from the weakly-supervised approach, as it is trained in a fully supervised way (instead of weakly-supervised) with coarse annotation and pixel-level ground truth in pairs. Additionally, we note several existing mask refinement methods [21,5] which refine initial coarse masks into more accurate segmentation results; however, they are usually optimized for a particular dataset. In comparison, our method applies the prototype learning paradigm [17,7,14,24] to enhance generalization across different data sets and lesion types. Good generalization is the key to putting the coarse annotation refinement algorithm into practice. For example, we can train the coarse annotation refinement network on a large-scale dataset for once and reuse the trained model on other datasets with less fine-tuning or tweakings.

Particularly, our prototype learning averages the features from the coarse mask region to form a lesion prototype and averages the background features to create a background prototype. A pixel is classified to the lesion class if its corresponding feature vector is more similar to the lesion prototype. Since our prototypical approach generates image-specific prototype to adaptively describe the image itself, it is less sensitive to the intra-class variance and high distribution shifts from different datasets or unseen classes. However, averaging features uniformly may be problematic when the lesion is considerably smaller than the coarse mask, as the resultant lesion prototype becomes dominated by background features. We alleviate this issue by a superpixel-guided prototype weighing module. The module first divides the coarse mask into several superpixels[7] and the prototype for each superpixel is obtained. Each prototype’s dis-similarity with the background prototype is then calculated as a weighting factor. The final lesion prototype is the weighted combination of these superpixel-guided sub-prototypes.

Contributions. (1) To the best of our knowledge, our method is the first prototypical approach for the coarse retinal lesion annotation refinement problem. (2) We present a prototype weighing module to solve the problem of the actual lesion being overly small. (3) Experiments demonstrate that the proposed method substantially improved the initial coarse annotation and outperformed non-prototypical mask refinement baselines. It also confirms the superiority of the prototype weighing module in both cross-dataset and cross-class settings.

2 Methods

This section details the proposed coarse annotation refinement method with the overall structure shown in Fig. 1. We assume that there exists a set of image patches and the associated coarse lesion annotations, and our algorithm will convert them into the corresponding high-quality pixel-level annotations.

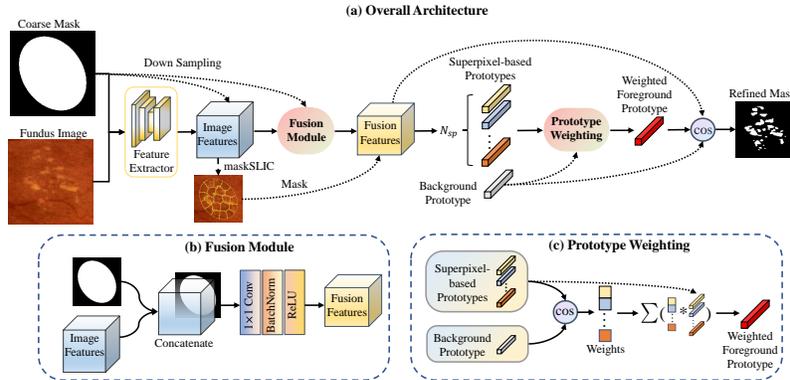


Fig. 1. Framework of our prototype-based coarse annotation refinement network.

2.1 Annotation Refinement via Prototype Learning

Feature extraction The input to the network is the concatenation of the image patch $I \in \mathbb{R}^{H \times W \times 3}$ and its corresponding coarse lesion annotation $M \in \mathbb{R}^{H \times W \times 1}$. We use a modified U-Net backbone to extract its feature map $F \in \mathbb{R}^{H' \times W' \times C}$. Following [14], we remove the last two upsampling blocks in the U-Net to speed up the calculation. As a result, the resolution of the feature map is 1/4 of the original input. We concatenate the feature map with the down-sampled coarse mask $M' \in \mathbb{R}^{H' \times W' \times 1}$ in the feature channel dimension to further incorporate the coarse annotation prior. To get the final fusion feature map $F' \in \mathbb{R}^{H' \times W' \times C'}$, we adopt a simple 1-layer network with architecture: 1×1 Conv2d+BatchNorm2d+ReLU.

Coarse Prototype Extraction Given the fused feature map, we want to learn representative and well-separated prototype vectors for the lesion region and the background based on the prototypical network. In previous research [17,14,22], the prototypical network condenses the masked object features in an image into a single or few prototypes. A relative simple coarse foreground lesion prototype can be calculated by mask average pooling, as follows:

$$p_{fg} = \frac{\sum_{(x,y)} F'(x,y) \mathbb{1}[M'(x,y) = 1]}{\sum_{(x,y)} \mathbb{1}[M'(x,y) = 1]}, \quad (1)$$

where (x, y) indexes the spatial locations and $\mathbb{1}(\bullet)$ is an indicator function. In addition, the background prototype is computed by

$$p_{bg} = \frac{\sum_{(x,y)} F'(x, y) \mathbb{1}[M'(x, y) = 0]}{\sum_{(x,y)} \mathbb{1}[M'(x, y) = 0]}, \quad (2)$$

where $p_{fg}, p_{bg} \in \mathbb{R}^{C'}$.

Coarse Annotation Refinement Refinement is done using a non-parametric metric learning method [17]. For each pixel at location (x, y) of the final fusion feature map F' , we calculate the distance between its feature vector and the derived prototypes $\mathcal{P} = \{p_{bg}, p_{fg}\}$. Then, we apply the softmax operation over the distances to get the probability map $P_c \in \mathbb{R}^{H' \times W' \times 1}$ and $c \in \{bg, fg\}$. Formally, we have:

$$P_c(x, y) = \frac{\exp(-\alpha \cdot d(F'(x, y), p_c))}{\sum_{p_j \in \mathcal{P}} \exp(-\alpha \cdot d(F'(x, y), p_j))}, \quad (3)$$

where α is the scaling factor fixed at 20.

We train our model end-to-end using the sum of dice loss \mathcal{L}_{dice} and binary cross-entropy loss \mathcal{L}_{bce} between the final probability map P_{fg} and the well-annotated ground truth mask M_{gt} . That is: $\mathcal{L}_{loss} = \mathcal{L}_{dice} + \mathcal{L}_{bce}$.

During testing, for each image patch I and its corresponding coarse lesion annotation M , we obtain a corresponding foreground probability map P_{fg} . When mapping P_{fg} back to the original image space (uncropped full image), some of them will overlap. For each pixel in the overlapping area, we choose the maximum probability of these probability maps as its value. In the end, thresholding is used to get the final refined mask.

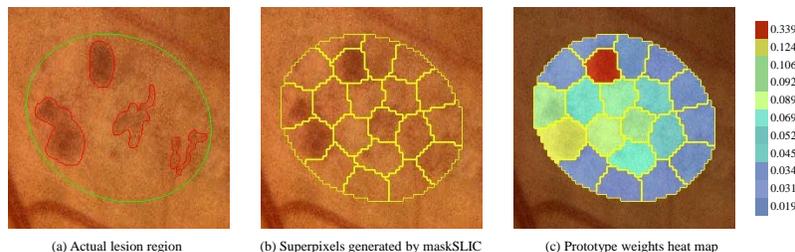


Fig. 2. Illustration of superpixel-guided prototype weighting. Boundaries in red, green, and yellow are actual lesion, coarse lesion, and superpixel regions respectively

2.2 Superpixel-guided prototype weighting

As shown in Fig. 2 when the actual lesion is relatively small compared to the coarse mask, the coarse foreground prototype defined by Eq. (1) cannot rep-

resent the actual lesion features. To reduce the impact of false-positive pixels within the coarse annotation, we divide the initial coarse region into several sub-regions according to their feature similarity. Concretely, we refer to maskSLIC [6] to aggregate the feature map within the masked region into multiple superpixel clusters. For each superpixel region S_i , we can obtain its corresponding superpixel-based sub-prototype g_i according to Eq. (1). We collect the extracted sub-prototypes and denote them as set $\mathcal{G} = \{g_i\}$, where $i \in 1, 2, \dots, N_{sp}$ (N_{sp} is the number of superpixels). We compute the cosine distance to measure the similarity between each g_i and p_{bg} :

$$d(g_i, p_{bg}) = 1 - \frac{g_i \cdot p_{bg}}{\|g_i\| \cdot \|p_{bg}\|}. \quad (4)$$

Intuitively, the prototypes dissimilar to the background prototype are more important parts for the final foreground prototype. Therefore, we can derive a weight coefficient for each prototype in set \mathcal{G} :

$$w_i = \frac{\exp(\beta \cdot d(g_i, p_{bg}))}{\sum_{g_j \in \mathcal{G}} \exp(\beta \cdot d(g_j, p_{bg}))}, \quad (5)$$

where β is the scaling factor fixed at 10. The final foreground prototype is then given by

$$p_{\text{weighted}} = \sum_{g_i \in \mathcal{G}} w_i \cdot g_i. \quad (6)$$

As shown later, our proposed superpixel weighted prototype p_{weighted} is a more representative foreground prototype that performed better in various experiments.

3 Experiments and Results

3.1 Experimental Setup

Coarse Annotation Generation There is no public available retinal lesion dataset with paired coarse annotation and pixel-level segmentation mask. To construct such paired dataset, we develop a simple coarse annotation generation method. Firstly, the coarse annotations are simulated from the well-annotated fine masks by applying the following chain of operations: smoothing, dilating, expanding, clustering the connected components using DBSCAN [2] and fitting ellipses to each cluster. Secondly, the fundus image is cropped around each ellipse in the corresponding coarse annotation. Finally, these cropped image patch and coarse annotation pairs are resized to fixed dimensions $H \times W$ for subsequent model training and testing.

Datasets and Evaluation Metrics. We evaluate the proposed methods on publicly IDRiD and DDR datasets, and our real-world private dataset. IDRiD contains 81 fundus images (54 training images, 27 testing images) with pixel-level

annotations for hard exudates (EX), hemorrhages (HE), microaneurysms (MA), and soft exudates (SE). Similarly, the testing part of DDR contains 225 fundus images with pixel-level annotations for EX, HE, MA, and SE. Our real-world private testing dataset collects 211 fundus images with pixel-level annotations for drusen (Drus) and pre-retinal hemorrhages (Prh) labeled by two experienced ophthalmologists. To train our annotation refinement network, we collect 32985 training patch pairs (EX:9957, HE:7752, MA:14387, SE:889) from the IDRiD training images using our coarse annotation generation algorithm. We also apply the mask generation algorithm to generate coarse mask for each testing image. To compare different refinement methods, we calculate the Intersection over Union (IoU) between the refined annotation and ground-truth mask.

Baselines. We implement several non-prototypical mask refinement baseline models, taking in an image patch and a coarse mask as the input. In detail, we choose three widely used feature extraction backbones, Res18 [3], HRNet18 [23], and U-Net [13] attached with the coarse mask fusion module to perform feature extraction. The feature extraction process is identical to the one described in Sec.2.1 except for the feature backbone. After that, we attach a 1x1 Conv2d layer as a binary classifier to obtain a refined segmentation score map.

Implementation Details. For the prototypical methods, we set the superpixel number $N_{sp} = 20$. All training patch pairs are resized to 256×256 and augmented by RandomShiftScaleRotate, RandomBlur, and RandomBrightnessContrast. All models are implemented by PyTorch and trained from scratch using Adam optimizer with a batch size of 64 for 120 epochs. The initial learning rate is 10^{-4} and reduces according to ReduceLROnPlateau strategy.

3.2 Results

Same-Dataset Experiments We train the proposed coarse annotation refinement network using all four lesion types on IDRiD and evaluate the performance on the IDRiD testing set. As shown in Table 1, “Initial Coarse” denotes the IoU of actual lesion region versus coarse annotation region. Our prototypical method improves the initial coarse mask considerably. It also consistently obtains better refinement performance than the non-prototypical baselines on all four lesion classes in terms of IoU score, with average IoU score improving by more than 5.2%. This experiment demonstrates our advantages when training and testing images are from the same dataset.

Cross-Dataset and Cross-Class Experiments We directly evaluate the performance on the DDR testing set and our real-world private dataset using models trained on the IDRiD dataset without further fine-tuning. As shown in Table 2, our method exceeds the U-Net baseline by 4.3% on both DDR and private datasets. For DDR, our superpixel weighted prototype performs better for all

Table 1. The image-level average IoU (%) and its standard deviations (%) of on IDRiD. "w/ superpixel" means with superpixel-guided prototype weighing.

| Methods | MA | SE | EX | HE | Average |
|----------------|-------------------|-------------------|-------------------|--------------------|-------------------|
| Initial Coarse | 9.6 (2.6) | 49.3 (10.4) | 15.0 (5.2) | 33.2 (8.6) | 26.8 (6.7) |
| Res18 [3] | 73.9 (7.3) | 68.4 (14.2) | 54.1 (9.4) | 62.6 (10.9) | 64.7 (10.5) |
| HRNet18 [23] | 79.1 (7.5) | 78.1 (5.3) | 56.8 (9.1) | 64.8 (11.8) | 69.7 (8.4) |
| U-Net [13] | 77.6 (6.8) | 75.6 (12.6) | 58.9 (8.9) | 67.9 (11.4) | 70.0 (9.9) |
| Our methods | 84.2 (6.2) | 80.7 (9.2) | 65.3 (8.1) | 69.9 (13.9) | 75.0 (9.4) |
| w/ superpixel | 84.1 (6.2) | 79.6 (9.9) | 65.9 (8.7) | 71.1 (11.2) | 75.2 (9.0) |

lesion types compared to the non-weighted prototype. Similarly, the weighted prototype is notably better than the non-weighted one on the private dataset, especially for the class Prh (56.3% \rightarrow 62.6%). Overall, we see a general trend that our model can generalize well to new datasets or unseen classes.

Table 2. The image-level average IoU (%) and its standard deviations (%) on DDR and our real-world private dataset having 2 unseen classes.

| Methods | DDR | | | | | Private | | |
|----------------|--------------------|--------------------|--------------------|--------------------|--------------------|--------------------|--------------------|--------------------|
| | MA | SE | EX | HE | Average | Drus | Prh | Average |
| Initial Coarse | 6.9 (3.4) | 32.3 (10.8) | 14.8 (8.8) | 23.1 (11.4) | 19.3 (8.6) | 33.6 (19.1) | 49.7 (13.0) | 41.6 (16.0) |
| Res18 [3] | 56.3 (12.4) | 67.6 (13.6) | 52.3 (13.2) | 54.4 (13.3) | 57.6 (13.1) | 43.9 (16.6) | 52.9 (15.2) | 48.4 (15.9) |
| HRNet18 [23] | 65.1 (13.5) | 68.2 (12.6) | 54.2 (13.9) | 58.5 (12.9) | 61.5 (13.3) | 44.2 (21.1) | 57.5 (16.4) | 50.9 (18.8) |
| U-Net [13] | 60.9 (13.2) | 70.2 (15.0) | 55.4 (12.8) | 59.8 (12.7) | 61.6 (13.4) | 45.2 (20.9) | 57.9 (16.4) | 51.6 (18.6) |
| Our methods | 68.8 (12.0) | 71.2 (16.9) | 58.4 (12.7) | 61.3 (15.7) | 64.9 (14.3) | 47.9 (23.4) | 56.3 (20.7) | 52.1 (22.1) |
| w/ superpixel | 69.8 (11.9) | 72.0 (17.5) | 58.9 (12.4) | 62.9 (14.5) | 65.9 (14.1) | 49.1 (23.0) | 62.6 (15.3) | 55.9 (19.1) |

Coarse Mask Reduction Factors Since ophthalmologists tend to draw a single rough ellipse to cover several unconnected lesion regions, we simulate the process by setting different reduction factors to the DBSCAN clustering algorithm. Actually, the number of the generated ellipses is the number of connected lesion regions divided by the reduction factor. In other words, with a higher reduction factor, the generated coarse mask will be more coarse. As shown in Fig. 3, although the refinement performance of all methods degrades as the reduction factor ranges from 1.0 to 2.0, our prototypical method has less degradation compared to the U-Net baseline. It implies our method is more robust against coarser annotations.

Visual results Fig. 4 presents some visualization of refinement results. Despite the vast variation in lesion scales, colors, and low contrast to surrounding regions, the first three rows show our proposed superpixel weighted prototype approach generates the most accurate lesion boundary. The last row shows a failure case when the coarse mask contains two distinct lesion classes, EX and Drus, at the same time.

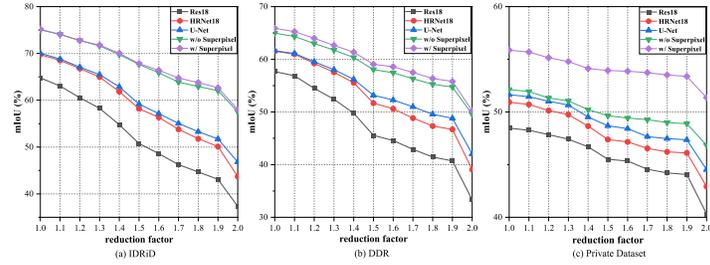


Fig. 3. Refinement performance under different reduction factors.

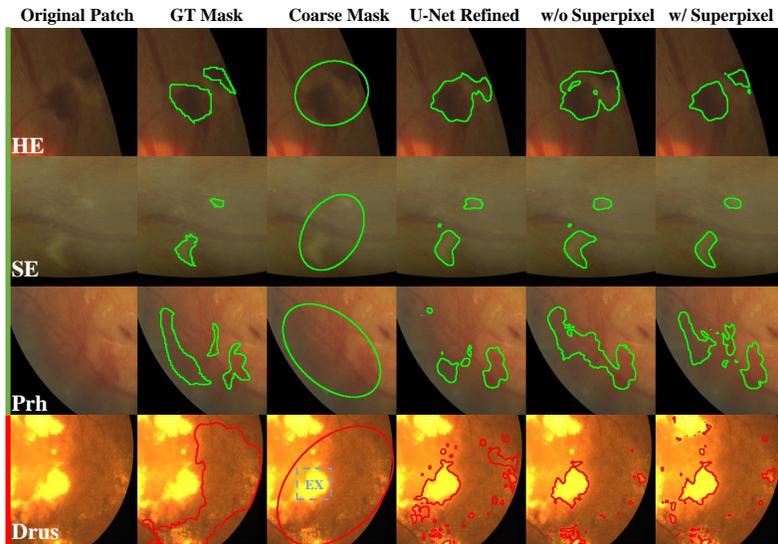


Fig. 4. Visualization of refinement results on four categories of lesions.

4 Conclusions

This paper proposes a novel prototype-based network to convert a coarse annotation into a pixel-level segmentation mask. The proposed network first extracts the lesion and background prototypes and labels the image pixel as the lesion class if its feature is more similar to the lesion prototype. A superpixel-guided prototype weighing module is then proposed to tackle the issue of the actual lesion being overly small compared to the coarse mask. On the IDRiD dataset, our model outperformed non-prototypical baselines by a large margin. Extensive experiments on DDR and our real-world private dataset also demonstrate the proposed model enjoys better generalizability to new datasets and some unseen lesion classes.

References

1. Chu, T., Li, X., Vo, H.V., Summers, R.M., Sizikova, E.: Improving weakly supervised lesion segmentation using multi-task learning. In: *Medical Imaging with Deep Learning*. pp. 60–73. PMLR (2021)
2. Ester, M., Kriegel, H.P., Sander, J., Xu, X., et al.: A density-based algorithm for discovering clusters in large spatial databases with noise. In: *kdd*. vol. 96, pp. 226–231 (1996)
3. 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)
4. Huang, S., Li, J., Xiao, Y., Shen, N., Xu, T.: Rtnet: Relation transformer network for diabetic retinopathy multi-lesion segmentation. *IEEE Transactions on Medical Imaging* (2022)
5. Huang, Y., Lin, L., Li, M., Wu, J., Cheng, P., Wang, K., Yuan, J., Tang, X.: Automated hemorrhage detection from coarsely annotated fundus images in diabetic retinopathy. In: *2020 IEEE 17th International Symposium on Biomedical Imaging (ISBI)*. pp. 1369–1372. IEEE (2020)
6. Irving, B.: maskslc: regional superpixel generation with application to local pathology characterisation in medical images. *arXiv preprint arXiv:1606.09518* (2016)
7. Li, G., Jampani, V., Sevilla-Lara, L., Sun, D., Kim, J., Kim, J.: Adaptive prototype learning and allocation for few-shot segmentation. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. pp. 8334–8343 (2021)
8. Li, T., Gao, Y., Wang, K., Guo, S., Liu, H., Kang, H.: Diagnostic assessment of deep learning algorithms for diabetic retinopathy screening. *Information Sciences* **501**, 511–522 (2019)
9. Liu, Q., Liu, H., Liang, Y.: M2mrf: Many-to-many reassembly of features for tiny lesion segmentation in fundus images. *arXiv preprint arXiv:2111.00193* (2021)
10. Liu, X., Yuan, Q., Gao, Y., He, K., Wang, S., Tang, X., Tang, J., Shen, D.: Weakly supervised segmentation of covid19 infection with scribble annotation on ct images. *Pattern recognition* **122**, 108341 (2022)
11. Playout, C., Duval, R., Cheriet, F.: A novel weakly supervised multitask architecture for retinal lesions segmentation on fundus images. *IEEE transactions on medical imaging* **38**(10), 2434–2444 (2019)
12. Porwal, P., Pachade, S., Kokare, M., Deshmukh, G., Son, J., Bae, W., Liu, L., Wang, J., Liu, X., Gao, L., et al.: Idrid: Diabetic retinopathy–segmentation and grading challenge. *Medical image analysis* **59**, 101561 (2020)
13. Ronneberger, O., Fischer, P., Brox, T.: U-net: Convolutional networks for biomedical image segmentation. In: *International Conference on Medical image computing and computer-assisted intervention*. pp. 234–241. Springer (2015)
14. Tang, H., Liu, X., Sun, S., Yan, X., Xie, X.: Recurrent mask refinement for few-shot medical image segmentation. In: *Proceedings of the IEEE/CVF International Conference on Computer Vision*. pp. 3918–3928 (2021)
15. Tang, Y., Cai, J., Yan, K., Huang, L., Xie, G., Xiao, J., Lu, J., Lin, G., Lu, L.: Weakly-supervised universal lesion segmentation with regional level set loss. In: *International Conference on Medical Image Computing and Computer-Assisted Intervention*. pp. 515–525. Springer (2021)
16. Wang, J., Xia, B.: Bounding box tightness prior for weakly supervised image segmentation. In: *International Conference on Medical Image Computing and Computer-Assisted Intervention*. pp. 526–536. Springer (2021)

17. Wang, K., Liew, J.H., Zou, Y., Zhou, D., Feng, J.: Panet: Few-shot image semantic segmentation with prototype alignment. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 9197–9206 (2019)
18. Wei, Q., Li, X., Yu, W., Zhang, X., Zhang, Y., Hu, B., Mo, B., Gong, D., Chen, N., Ding, D., et al.: Learn to segment retinal lesions and beyond. In: 2020 25th International Conference on Pattern Recognition (ICPR). pp. 7403–7410. IEEE (2021)
19. Yan, Z., Han, X., Wang, C., Qiu, Y., Xiong, Z., Cui, S.: Learning mutually local-global u-nets for high-resolution retinal lesion segmentation in fundus images. In: 2019 IEEE 16th International Symposium on Biomedical Imaging (ISBI 2019). pp. 597–600. IEEE (2019)
20. Yang, L., Zhang, Y., Zhao, Z., Zheng, H., Liang, P., Ying, M.T., Ahuja, A.T., Chen, D.Z.: Boxnet: Deep learning based biomedical image segmentation using boxes only annotation. arXiv preprint arXiv:1806.00593 (2018)
21. Yang, Y., Wang, Z., Liu, J., Cheng, K.T., Yang, X.: Label refinement with an iterative generative adversarial network for boosting retinal vessel segmentation. arXiv preprint arXiv:1912.02589 (2019)
22. Yu, Q., Dang, K., Tajbakhsh, N., Terzopoulos, D., Ding, X.: A location-sensitive local prototype network for few-shot medical image segmentation. In: 2021 IEEE 18th International Symposium on Biomedical Imaging (ISBI). pp. 262–266. IEEE (2021)
23. Yuan, Y., Chen, X., Wang, J.: Object-contextual representations for semantic segmentation (2020)
24. Zhang, G., Lu, X., Tan, J., Li, J., Zhang, Z., Li, Q., Hu, X.: Refinemask: Towards high-quality instance segmentation with fine-grained features. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 6861–6869 (2021)