

The 2nd LID Challenge (Weakly Supervised Object Localization)

Zhendong Wang, Zhenyuan Chen, Chen Gong Nanjing University Of Science and Technology LEAP Group@PCA Lab



Outline







3 Experimental Results



Problem Analysis







- Weakly-Supervised Object Localization: localize the objects in an image only with imagelevel annotations.
- Only using image-level labels is too weak, so we try to utilize the information of raw image to promote the edge and other details.
- We designed a network composed of two auto-encoder parts.

Model Description





- We designed a network with two auto-encoders.
- In the first part, we train a classifier with global average pooling under the supervision of image-level annotations. Then we use the binary images generated by CAMs as pseudo pixel-level annotations.
- In the second part, we expect to recover the raw image from binary image in order to get the refined binary image for the next iteration.

Model Description



Our loss function:

$$L = \sum_{k=1}^{K} a_k l^k + \frac{1}{2} \sum_{i=1}^{c} (\hat{y}_i - y)^2 + l_{out}$$

$$l_{out} = -\sum_{(r,c)} \left[I(r,c) \log(P(r,c)) + (1 - I(r,c)) \log(1 - P(r,c)) \right]$$

$$l^k = l_{bce}^k + l_{ssim}^k$$

$$l_{bce} = -\sum_{(r,c)} \left[G(r,c) \log(S(r,c)) + (1 - G(r,c)) \log(1 - S(r,c)) \right]$$

$$l_{ssim} = 1 - \frac{(2u_x u_y + 0.01^2)(2\sigma_{xy} + 0.03^2)}{(u_x^2 + u_x^2 + 0.01^2)(\sigma_x^2 + \sigma_y^2 + 0.03^2)}$$

- 1. The loss l^k between binary image prediction and pseudo pixel-level annotations.
- 2. Mean square loss of class prediction.
- 3. Cross entropy loss I_{out} between the final output image and the raw image.



Rank 🜲	Participant team 🝦	Peak_loU 🌲	Peak_Threshold ᅌ	Last submission at 👙
1	VL-task3	0.63	24.00	2 days ago
2	BJTU-Mepro-MIC	0.62	35.00	2 days ago
3	LEAP Group@PCA Lab	0.61	7.00	2 days ago

- Our model achieved 61% Peak_IoU in test dataset.
- Because of the wrong choice of one param of output image function, the Peak_Threshold is only 7.
- We corrected this mistake and improved the Peak_Threshold.





Above output images show that our method works well on images with a single object or overlapping objects.

Experiment Results







- The three images in left panel demonstrate that our model can localize the objects with small local complex edge structure clearly.
- The three images in right panel show that our method is also applicable to the images containing multiple instances belonging to one category.

Failure cases:



We found our method does not work well on situations as shown above. One is with special strips, and the other is with some interfering information.

Future Improvements

A LOOTON SCIENCE

Our plan:

1. Bringing Res2Net structure to downsampling layers. Providing information of rich scales by Integrating the channels.





2. Add Strip Pooling in our encoderdecoder model, as Strip Pooling can help the network better exploit long-range dependencies.

Future Improvements





We carried out some preliminary experiments. The results show that our improvements promote the effects on some challenging images.

THANKS FOR LISTENING