arXiv:2008.01786v1 [cs.CV] 4 Aug 2020

# Entropy Guided Adversarial Model for Weakly Supervised Object Localization

Sabrina Narimene Benassou[a], Wuzhen Shi[b], Feng Jiang[a],*

[a]*92 Xidazhi Street, Nangang District, Harbin City, Heilongjiang Province, Harbin Institute of Technology*
[b]*College of Electronics and Information Engineering, Shenzhen University No. 3688, Nanhai Avenue, Nanshan District, Shenzhen, China*

## Abstract

Weakly Supervised Object Localization is challenging because of the lack of bounding box annotations. Previous works tend to generate a class activation map i.e CAM to localize the object. Unfortunately, the network activates only the features that discriminate the object and does not activate the whole object. Some methods tend to remove some parts of the object to force the CNN to detect other features, whereas, others change the network structure to generate multiple CAMs from different levels of the model. In this present article, we propose to take advantage of the generalization ability of the network and train the model using clean examples and adversarial examples to localize the whole object. Adversarial examples are typically used to train robust models and are images where a perturbation is added. To get a good classification accuracy, the CNN trained with adversarial examples is forced to detect more features that discriminate the object. We futher propose to apply the shannon entropy on the CAMs generated by the network to guide it during training. Our method does not erase any part of the image neither does it change the network architecure and extensive experiments show that our Entropy Guided Adversarial model (EGA model) improved performance on state of the arts benchmarks for both

---

localization and classification accuracy.

## 1. Introduction

Weakly supervised learning has gained a lot of popularity during these past few years, especially since Zhou et al [1] propose the use of Class Activation Map (CAM) to localize objects without bounding boxes annotations. Since that, CAM was extensively used for object localization [2, 3, 4, 5, 6, 7, 8], object detection [9, 10, 11, 12, 13, 14], image segmentation [15, 16, 17], etc. However, not the whole object is highlighted on the CAM, because the network learns the features that discriminate the most the object. Some works have been proposed to deal with this probelm, we can classify them into two approaches. The first approach [2, 3, 4, 8] consists of hiding a part or some parts of the image during training to force the CNN to detect the full object, nevertheless this approach has a drawback, when the most discriminative part is removed the CNN tends to learn regions of the image that does not belong to the object (as water or tree branches in CUB dataset) because they appear frequently in the training samples [4, 18]. Furthermore, removing some parts of the image results in information loss for the network and hence decreases its recognition ability [8]. The second approach [5, 6, 7] consists of activating different parts of the object by generating multiple CAMs from different levels of the network. These methods however, require a modification in the network structure by plugging some layers or some blocks to the network, which is not always intuitive to design or to generalize for different network structures. In this paper, we propose to take advantage of the generalization ability of the network and propose to use Adversarial Learning (AL) [19, 20] and entropy [21, 22] to tackle the limitations of the two approches. Our method does not modify the network backbone, which make the implementation easier, and because the whole images are used for training, there is no information loss for the model.
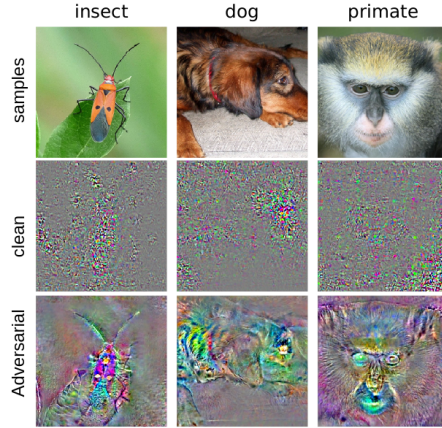
Figure 1: The difference between the features that affect the classification prediction when a model is trained with clean examples and a model trained with adversarial (noisy) examples [23]

Adversarial learning was originally made to prevent adversarial attacks by training the model with adversarial examples that are generated by introducing some perturbations to clean images. In [23], Tsipras et al. visualize the gradient of the loss on ImageNet dataset and have shown the features that affect the classification prediction of a CNN trained with clean examples and a CNN trained with adversarial (noisy) examples. As we can see in Fig. 1, the CNN trained adversarially, detect more relevant and clearer features such as edges and borders (line 3) than the standard one, that detects noisy features (line 2). This leads us to the following conclusion; to get a good classification accuracy, the CNN trained only with clean images activate the features that better discriminate the object e.g the head of a bird. But if we futher train the CNN with images where we add a small perturbation, the network is forced to look for more relevant features to better recognize the object e.g the body of the bird. By training a model with both clean and adversarial examples, the generated CAM will activate more discriminative pattern to recognize the object without erasing any part of the image or changing the network architecture. Furthermore, Goodfellow et al. [20] advice the use of a mixture of clean and

3

adversarial examples to train the CNN. Data augmentation such as rotation or translation will not occur naturally and treating adversarial examples as data augmentation will in fact regularize the network.

As adversarial learning detects more discriminative features of the object, some remaining pixels that belong to the object are still not activated by the CAM. To this end, we propose to further use the concept of entropy to guide the CNN during training. The entropy is a measure of uncertainty, it represents how a CNN is certain or not about its predictions [21, 9, 22], the CAM generated by the CNN only highlights the most discriminative features of the object. The pixels that constitute this part have a high prediction probability, i.e., low-entropy. And the other pixels, that are not highlighting by the CAM have a low prediction propability, i.e., high-entropy. By minimizing the entropy of the CAM, the CNN could extend the localization of the object features.

We can resume our work to the following contributions :

1. An entropy guided adversarial model (Dubbed EGA model) that uses both clean examples and adversarial examples to activate more features of the object.

2. Introduce an entropy loss function to guide the CNN to detect the pixels that belong to the object.

3. Extensive expriments shown that our EGA model obtained state-of-theart performance in weakly supervised object localization.

## 2. Related work

### 2.1. Weakly Supervised Object Localization

Weakly supervised object localization recieved a lot of attention since [1], where they used a Global Average Pooling (GAP) layer to generate a Class Activation Map that localizes the object with only label annotation. This map highlights only the most discriminative part of the object, resulting in a tight bounding box. Many works attempt to create new methods to solve this problem. We can devide the proposed solutions into two classes, the first class

remove a part or some parts of the image to force the network to detect more features. As in [2], where they proposed to use two branches, the first branch generates a map with the discriminative part highlighted, and then removes the pixels of this part, to give it as input to the second branch of the network, the second branch generates another map where other features of the object are highlighted. The two generated maps are then fused by taking the maximum value. Another method that uses a hiding process is [3], in which the images are divided into patches, and then for each patch a probability is assigned. At each epoch, some patches are hidden and given as input to the CNN. The CNN hence learns to detect the whole object. During testing, the images are given to the network without removing any patches. [4] uses either an attention mechanism or a dropout mask to help the network to detect the whole object. [8] removes the discriminative region and replace it with a patch from another image, hence there is no non-informative pixels in the image and the network could generalize well for both classification and localization. Some other methods do not remove any part of the image and modify the network achitecture to generate different maps from different hierarchies. In [5], an attention map generated from high level features map, is used to guide the network to distinguish between the forground and background pixels of the map generated from low level features map. [6] combines two child classes labels to create a parent class and train the CNN with hierarchical class labels to detect common visual patterns and [7] generates many CAMs from low and high features map and combine them using a polynomial function.

For the first approach, when some parts of the object are removed, non discriminative features are activated by the network because they appear too frequently in the dataset. The second approach moreover is not intuitive to apply for complex network architectures as we have to plug some blocks to generate multiple CAMs from different levels of feature maps. In this article, we propose to employ adversarial examples as data augmentation to detect more discriminative features. Minimizing the entropy loss on the CAM will further guide the model during training.

5

## 2.2. Adversarial learning

Adversarial learning consists of constructing robust models by training the classifier with adversarial examples [20]. Some works used it either to improve the robustness of the model, or to solve other problems; as Madry et al. [19] where they constructed a robust classifier using a min-max formulation to protect their network against adversarial attacks. Tsipras et al. [23] relates the benefit of adversarial learning and show that robust classifier learns different features than standard one. [24] trains a network with both clean examples and adversarial examples to improve the classification accuracy and achieved state of the art accuracy on ImageNet dataset. [25] employs adversarial learning to improve classification in semi-supervised learning. In [26], an adversarial dropout mask is selected using adversarial learning and applied to the network to improve the classification accuracy for both supervised and semi-supervised learning. [27] proposes a fast training for adversarial training by updating the network parameters and creating an adversarial example in one backward pass. [28] applies adversarial learning and adversarial dropout to learn discriminative features for unsupervised domain adaptation. In this article, we take advantage of adversarial attacks and use adversarial examples as data augmentation to improve performance of weakly supervised object localization problems.

## 2.3. Entropy

In information theory, entropy is the measure of uncertainty, it tells how a network is uncertain about its prediction. When the prediction of the network is of high porbability, the entropy value is low, and when the prediction is low, the amount of entropy is high. [9] uses entropy for weakly supervised object detection as an optimization method associated with multiple instance learning to select the object proposals that belong to the object. [21] shows state of the art performance in Domain Adaptation in Semantic Segmentation using an entropy minimization loss on the target dataset. [22] considers the entropy as a weighting coefficient to improve weakly supervised learning on Pascal VOC
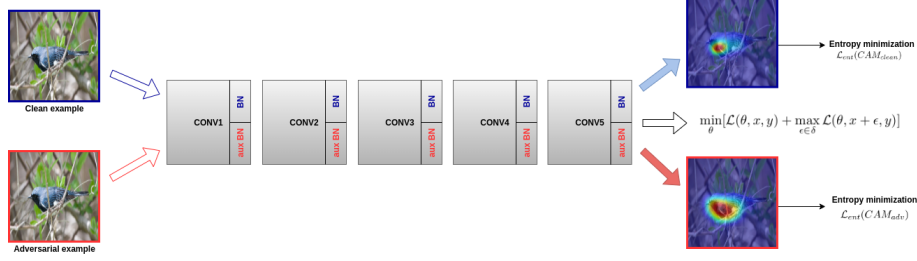
Figure 2: The clean example and adversarial example are fed through the network. The clean example passes through the main batch normalization and the adversarial learning through the auxiliary batch normalization. After the forward pass, we get $CAM_{clean}$ and $CAM_{adv}$ and calculate the shannon entropy loss.

dataset [29]. Entropy was also used for semi-supervised domain adaptation [30] and semi-supervised learning [31].

## 3. Method

In this section, we introduce our Entropy Guided Adversarial model i.e EGD model. In section 3.1, we review the baseline for WSOL. In 3.2, we briefly present adversarial learning and how to apply it for WSOL problem. In 3.3, we discuss about the shannon entropy application on CAM, and finally in section 3.3, we present the loss function used by our network. Our model is illustrated in Fig. 2.

### 3.1. Revising Class Activation Map

Weakly Supervised Object Localization aims at detecting objects without bounding box annotations. One method widely used to localize objects using only label annotations is [1]. In [1], they propose to add to a network, composed of convolutional layers a global average pooling layer before the softmax layer. We denote the last convolutional layer as $f_k(h, w)$. The GAP layer is applied to the last convolutional layer and output a vector where each unit is the average of each feature map :
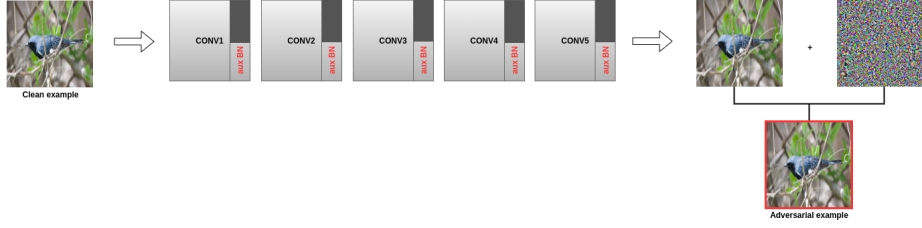
$$F_k = \sum_{(h,w)} f_k(h, w) \tag{1}$$

Figure 3: The clean example is given as input to the network to create a perturbation $\epsilon$ that maximizes the loss function. We drop the main batch norm and use the auxiliary batch norm to generate the adversarial example.

The weights of the generated vector are then multiplied by the features map of the last convolutional layer and are given as input to the softmax layer as:

$$S_c = \sum_k a_k^c \sum_{(h,w)} f_k(h, w) \qquad (2)$$

where $a_k^c$ is the weight of the corresponding class c which indicates the importance of $F_k$ for class c.

To generate a map that indicates the importance of each pixel for the corresponding class, we summed up the weighted features map as:

$$CAM_c(h, w) = \sum_k a_k^c f_k(h, w) \qquad (3)$$

where $CAM_c(h, w)$ is the class activation map generated for class c. The CAM activates the features that discriminate the most the object.

*3.2. Adversarial learning*

To train a CNN, we use stocastic gradient descent method to minimize an objective function according to the parameters $\theta$ of the network:

$$\min_\theta \mathcal{L}(\theta, x, y) \qquad (4)$$

where $x \in \mathcal{X}$ is an input sample associated with its ground truth label $y \in \mathcal{Y}$ and $\mathcal{L}(., ., .)$ is the objective function i.e loss function.

To construct a model that is robust to adversarial attacks, we train the classifier with adversarial examples by adding a small perturbation to the image
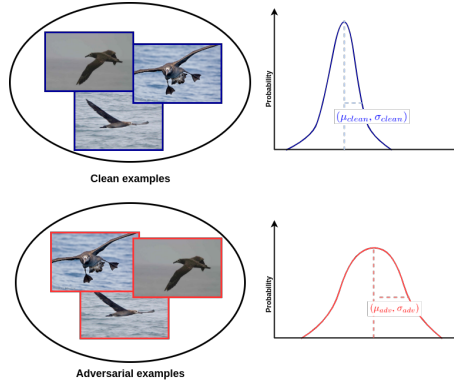
Figure 4: The difference between the distribution of clean examples and adversarial examples.

to fool the network [20, 19]. The resulted sample should be close to the original image, and maximize the loss function instead of minimizing it as we usually do for standard training, the objective function hence becomes :

$$\min_{\theta}[\max_{\epsilon \in \delta} \mathcal{L}(\theta, x + \epsilon, y)] \tag{5}$$

where $\epsilon$ is the adversarial perturbation i.e noise and $\delta$ is the allowable set of perturbations which ensure that adversarial image is close to the clean one.

Goodfellow et al [20] advice the use of a mixture of clean examples and adversarial examples to train the model. Adversarial examples act as data augmentation and associated with clean examples regularize the network. As suggested in Goodfellow et al work the objective function is:

$$\min_{\theta}[\mathcal{L}(\theta, x, y) + \max_{\epsilon \in \delta} \mathcal{L}(\theta, x + \epsilon, y)] \tag{6}$$

Clean samples and adversarial samples can be considered as two different datasets, because they have two different disctributions, hence the network could not generalize well on clean data during the inference. This is due to the distribution mismatch between the two type of datasets [24]. As shown in Fig. 4, the mean and the standard deviation of the ditribution of clean data is different from the one of adversarial data.

To solve this problem, we follow the work of [24] and add an auxiliary batch normalization to the network. Batch normalization (BN) is a technique used to

9

normalize the input data [32]. As we have two datasets distributed differently, we use one batch norm for each type of dataset i.e we feed each sample to its corresponding batch norm. Specifically, for each clean mini batch $x_{clean}$, we generate its corresponding adversarial mini-batch $x_{adv}$ using the auxiliary batch norm layer. The process of generating the adversarial example is shown in Fig. 3. We feed a clean sample to the network to calculate the perturbation that maximizes the loss function using Eq (5). The network parameters are not updating during this process. We then pass the clean mini batch and adversarial mini batch through the network to compute the loss, the clean mini batch pass through the main batch normalization and the adversarial mini batch pass through the auxiliary batch normalization. After the forward pass, the model generates two maps; $CAM_{clean}$ and $CAM_{adv}$ for clean sample and adversarial sample respectively. After that, we calculate the loss function and update the network parameter. During inference, we test our model only on clean data, the auxialiary batch norm is dropped and only the main batch norm is used.

### 3.3. Entropy minimization

Unlike previous wroks [2, 1], where they generate the CAM during the test stage, we generate our CAMs after each forward pass to guide the CNN for localization during the training. Our adversarial model using the main and auxiliary batch norm layer generates two different CAMs; $CAM_{clean}$ and $CAM_{adv}$ for the clean example and adversarial example respectively. These CAMs activate different parts of the object but still do not activate the whole object. We propose to use the shannon entropy to remedy to this problem.

The Shannon entropy measure the amount of uncertainty [21], when it is applied to the CAM, a pixel with a low entropy means a high prediction probability and a pixel with a high entropy means a low prediction probability. The maps generated by the network highlight the most discriminative part i.e these pixels have a low entropy; the other pixels that are not highlighted by the maps but belongging to the object have a high entropy. By minimizing the entropy of the CAM, we force the CNN to activate the pixels that belong to the object

but not highlighted by the map

Given the generated CAMs, $CAM_{clean}$ and $CAM_{adv}$, we calculate the entropy loss for one $CAM$ as :

$$\mathcal{L}_{ent}(CAM) = - \sum_{(h,w)} P_{CAM^{(h,w)}} \log P_{CAM^{(h,w)}} \tag{7}$$

where $P_{CAM^{(h,w)}}$ represent the probability of the pixel $(h,w)$. $\mathcal{L}_{ent}(CAM)$ is the sum of all entropies pixels i.e we maximize the predictions on the pixels that are not activated by the CAM.

### 3.4. Loss function

After feeding the network with clean and adversarial sample, we generate $CAM_{clean}$ and $CAM_{adv}$, we then calculate the loss entropy from the two maps. We jointly optimize the adversarial learning loss function and entropy loss function to optimize the network parameters by summing Eq (6) and Eq (7). The final loss function is defined as below :

$$\min_{\theta}[\mathcal{L}(\theta, x, y) + \max_{\epsilon \in \delta} \mathcal{L}(\theta, x + \epsilon, y)] + \lambda_{CAM_{clean}}\mathcal{L}_{ent}(CAM_{clean}) + \lambda_{CAM_{adv}}\mathcal{L}_{ent}(CAM_{adv}) \tag{8}$$

where $\lambda_{CAM_{clean}}$ and $\lambda_{CAM_{adv}}$ is the weighting factor controlling the importance of $\mathcal{L}_{ent}(CAM_{clean})$ and $\mathcal{L}_{ent}(CAM_{adv})$, respectively.

Because of the perturbation on the adversarial example, $CAM_{adv}$ has more activated pixels than $CAM_{clean}$. To this end, we set $\lambda_{CAM_{clean}} > \lambda_{CAM_{adv}}$, this setting pushes the classifier to be more severe with clean examples by forcing the $CAM_{clean}$ to activate more pixels.

## 4. Experiments

### 4.1. Experimental Setup

*Datasets.* We evaluate our method on two commonly used datasets for WSOL i.e CUB-200-2011 [33] and ILSVRC [34, 35]. We further evaluation EGA model on OpenImages, a fresh new dataset proposed by [18] for WSOL. CUB-200-2011 consists of 200 bird species, this dataset contains 11,788 images with 5,994

images for training and 5,794 for testing. ILSVRC 2016, is a large scale dataset of 1000 classes, that comprise 1.2 million for training, and 5,000 images for the validation set that we use for testing. OpenImages has 100 classes, it contains 29 819 for training, 2 500 for validation, and 5 000 images for testing.

*Evaluation metrics.* For classification evaluation, we use the Top-1 classification accuracy, which indicates that a prediction is correct when the prediction of the model is equal to the ground-truth class. For localization, as CUB and ILSVRC datasets have bounding boxes annotations, we use three evaluation metrics : Top-1 localization accuracy [34], Correct Localization [36] (CorLoc) rate and MaxBoxAccV2 [18]. Top-1 localization accuracy counts a localization as correct when the predicted class is correct and the predicted bounding box has an Intersaction Over Union (overlap) with the ground truth bounding box greater than 0.5. Correct Localization (CorLoc) is the localization performance whether or not the predicted class is correct. MaxBoxAccV2 is a new metric proposed by [18], which is an improved version of CorLoc; we average the results of the IOU between the predicted box and the ground truth box accross 0.3, 0.5, 0.7. OpenImages provides pixel-wise annotations, hence we use pixel average precision (PxAP) metric [18] for localization, which measure the pixelwise precision and recall trade-off.

*Experimental details.* We use for training VGGnet [37] and GoogLeNet [38], We follow the same setting as [1], we remove the layers after conv5-3 (from pool5 to prob) of the VGG-16 network and the last inception block of GoogLeNet. We then add two convolutional layers with kernel size $3 \times 3$, stride 1, pad 1 with 1024 units, and a convolutional layer of size $1 \times 1$, stride 1 with 1000 units for ILSVRC, 200 units for CUB-200-2011 and 100 units for OpenImages. For training, input images are resized to $256 \times 256$, then randomly cropped to $224 \times 224$. Both backbone networks are fine-tuned on the pre-trained weights of ILSVRC. During inference, we resized images to $224 \times 224$ to find the whole objects and for classification for CUB and ILSVRC datasets, we average the scores from the softmax layer with 10 crops.

| Method | top1 cls-err | top1 loc-err |
|---|---|---|
| GoogLeNet-CAM [1] | 26.2 | 58.94 |
| GoogLeNet-SPG [5] | - | 53.36 |
| GoogLeNet-ADL [4] | **25.45** | **46.94** |
| GoogLeNet-DANet [6] | 28.8 | 50.55 |
| **GoogLeNet-EGA (ours)** | 27.89 | 54.26 |
| VGGnet-CAM [1] | 23.4 | 55.85 |
| VGGnet-ACoL [2] | 28.1 | 54.08 |
| VGGnet-SPG [5] | 24.5 | 51.07 |
| VGGnet-CutMix [8] | - | 47.47 |
| VGGnet-ADL [4] | 34.73 | 47.64 |
| VGGnet-DANet [6] | 24.6 | 47.48 |
| VGGnet-CCAM [7] | 26.8 | 49.93 |
| **VGGnet-EGA (ours)** | **21.87** | **40.84** |

Table 1: Comparison to the state-of-the-art performance on the CUB dataset.

## 4.2. Comparison with the state-of-the-arts

We compare our EGA model to the state of the arts on CUB, ILSVRC and the new proposed dataset OpenImages. The results are shown in Tab. 1, Tab. 2 and Tab. 3, respectively.

*CUB.* As shown in Tab. 1, with VGG model, our model outperform by far all the previous state of the arts in both classification and localization with 21.87% for classification and 40.84% for localization. We improved our baseline VGGnet-CAM by 1.53% and 15.01% for classification and localization respectively. We also surpass VGGnet-CCAM, the current state of the art for localization with a large margin of 9.09% for localization and margin of 4.93% for classification.

With a GoogLeNet backnone, our EGA model achieved 27.89% and 57.26% for classification and localization respectively. We did not achieve a new state of the art with GoogLeNet architecture but we surpass our baseline GoogLeNet-

| Method | top1 cls-err | top1 loc-err |
|---|---|---|
| GoogLeNet-Backprop [**?** ] | - | 61.31 |
| GoogLeNet-CAM [1] | 35.0 | 56.40 |
| GoogLeNet-HaS [3] | - | 54.53 |
| GoogLeNet-ACoL [2] | 29.0 | 53.28 |
| GoogLeNet-SPG [5] | - | 51.40 |
| GoogLeNet-ADL [4] | **27.17** | 51.29 |
| GoogLeNet-DANet [6] | 27.5 | 52.47 |
| **GoogLeNet-EGA (ours)** | 27.42 | **50.17** |
| VGGnet-Backprop [39] | - | 61.12 |
| VGGnet-CAM [1] | 33.4 | 57.20 |
| VGGnet-ACoL [2] | 32.5 | 54.17 |
| VGGnet-CutMix [8] | - | 56.45 |
| VGGnet-ADL [4] | 30.52 | 55.08 |
| NL-CCAM [7] | **27.7** | **49.83** |
| **VGGnet-EGA (ours)** | 29.36 | 52.69 |

Table 2: Comparison to the state-of-the-art performance on the ILSVRC validation set.

CAM for localization with 4.68%. We argue that GoogLeNet compared to VGG backbone is deeper and hence, needs a larger dataset as ILSVRC dataset to achieve good performance, as we add some perturbation to the samples, the network needs more data to improve its classification and localization accuracy.

*ILSVRC.* In Tab. 2, with VGG backbone, we achieved 29.36% and 52.69% for classification and localization respectively on ILSVRC dataset, we surpass the basline with a difference of 4.04% and 4.51%, we also suprpass all the state of the art method for both classification and localization, except for NL-CCAM which outperforms our method with a margin of 1.66% and 2.86% for classification and localization respectively.

With GoogLeNet architecture, EGA model achieved 27.42% and 50.17% for classification and localization respectivaly. As supposed earlier, our method
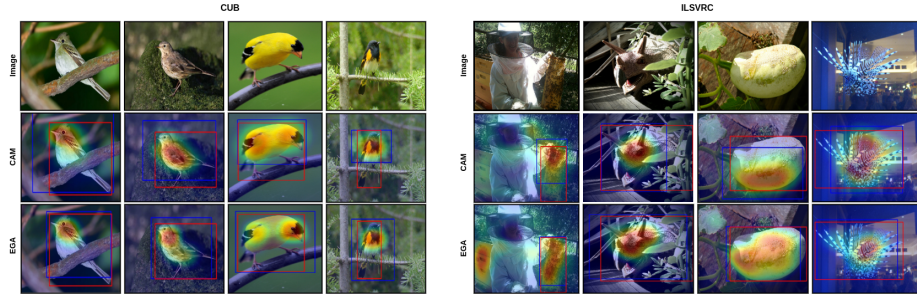
Figure 5: Comparison of our EGA method to CAM method. Our method activates more objects features than CAM and generates tighter bounding boxes. Ground-truth bounding boxes are in red and the predicted are in blue.
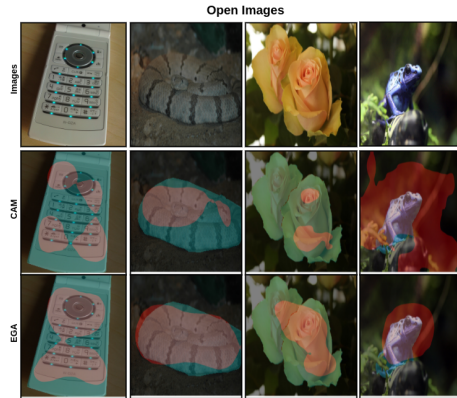


Figure 6: Comparison of our EGA method to CAM method on OpenImages dataset. The ground truth mask is in blue and the predicted mask is in red.

with GoogLeNet backbone with a larger dataset outperform all state of the art for localization, we improved our basline with a margin of 7.58% and 6.23% for classification and localization respectively. We also surpass the current state of the art ADL with a difference of 1.12% for localization with a good classification accuracy.

In Fig. 5, we compare the bounding boxes generated by CAM method [1] and bounding boxes generated by our EGA model. Our method activates more objects features than the baseline [1].

| Method | top1 cls-err | top1 loc-err |
|---|---|---|
| GoogLeNet-CAM [1] | 63.4 | **36.8** |
| GoogLeNet-HaS [3] | **31.6** | 41.9 |
| GoogLeNet-ACoL [2] | 59.3 | 42.8 |
| GoogLeNet-SPG [5] | 53.4 | 37.7 |
| GoogLeNet-ADL [4] | 53.4 | 43.2 |
| GoogLeNet-CutMix [8] | 46.9 | 37.5 |
| **GoogLeNet-EGA (ours)** | 33.4 | 37.35 |
| VGGnet-CAM [1] | 32.7 | 41.7 |
| VGGnet-HaS [3] | 40.0 | 41.9 |
| VGGnet-ACoL [2] | 31.8 | 45.7 |
| VGGnet-SPG [5] | **28.3** | 41.7 |
| VGGnet-ADL [4] | 33.9 | 41.3 |
| VGGnet-CutMix [8] | 31.9 | 41.9 |
| **VGGnet-EGA (ours)** | 30.0 | **38.21** |

Table 3: Comparison to the state-of-the-art performance on OpenImages dataset with PxAP metric.

*OpenImages.* As OpenImages is a new dataset proposed by [18], we compared our method to the results reported in [18]. In Tab. 3, with VGG backbone, our method achevies 30.0% for classification and 38.21% for localization, we still improve the baseline CAM with 2.7% for classification and 3.49% for localization. We further surpass ADL the current state of the art on OpenImages dataset with a margin of 3.9% for classification and a margin of 3.09% for localization. Hence we achieved a new state of the localization on OpenImages dataset with a small drop for classification. In Fig. 6, our EGA model detects more foreground features and less background features compared to VGGnet-CAM [1].

With GooleNet, our method achieves 33.4% for classification and 37.35% for localization, we have a slight decrease compared to the baseline CAM of 0.55% for localization, but we outperform HaS, ACoL, SPG, ADL and CutMix with a margin of 4.55%, 5.45%, 0.35%, 5.85%, and 0.15%, respectively.

| Method | ILSVRC | CUB |
|---|---|---|
| VGGnet-CAM [1] | 40.0 | 36.3 |
| VGGnet-HaS [3] | 39.4 | 36.6 |
| VGGnet-ACoL [2] | 42.6 | 42.6 |
| VGGnet-SPG [5] | 40.1 | 43.7 |
| VGGnet-ADL [4] | 40.2 | **33.7** |
| VGGnet-CutMix [8] | 40.6 | 37.7 |
| **VGGnet-EGA (ours)** | **38.22** | 36.07 |

Table 4: Comparison to the state-of-the-art performance on the ILSVRC, CUB datasets with MaxBoxAccV2 metric.

*MaxBoxAccV2.* we further evaluate our method with the new evaluation metric MaxBoxAccV2 proposed by [18]. The results are shown in Tab. 4. For CUB dataset with a VGG backbone, we still surpass all the sate of the art CAM, HaS, ACoL, SPG and CutMix with a difference of 0.23%, 0.53%, 6.53%, 7.63% and 1.63%, except for ADL where there is a difference of 2.37%. For ILSVRC dataset, we achieved a new state of the art for localization by surpassing all previous works with a localization accuracy of 38.22%.

*CorLoc.* We also evaluate our method with Correct Localization metric, a highly used evaluation metric in WSOL [36]. As shown in Tab. 5, EGA method improves CAM, HaS, ACoL and SPG methods by 6.17%, 3.63%, 1.87% and 0.14%. We also outperform the other state of the arts methods, except for NL-CCAM where we have a slight increase of 0.4%. NL-CCAM applied a Non-Local module to the VGG backbone and hence, this method changes the architecture of the VGG and does not use the original VGG backbone. Our method does not change the backbone of VGG or GoogleNet and we still have competetive results with NL-CCAM. When we compare CCAM method applied to the original basline i.e VGGnet-CCAM we outperform it by 1.25%.

| Method | ILSVRC |
|---|---|
| AlexNet-GAP [1] | 45.01 |
| AlexNet-HaS [3] | 41.25 |
| AlexNet-GAP-ensemble [1] | 42.98 |
| AlexNet-HaS-ensemble [3] | 39.67 |
| GoogLeNet-GAP [1] | 41.34 |
| GoogLeNet-HaS [3] | 38.8 |
| GoogLeNet-ACoL [2] | 37.04 |
| GoogLeNet-SPG [5] | 35.31 |
| VGGnet-CCAM [7] | 36.42 |
| NL-CCAM [7] | **34.77** |
| **GoogLeNet-EGA (ours)** | 35.17 |

Table 5: Comparison to the state-of-the-art performance on the ILSVRC validation set with CorLoc metric.

*4.3. Ablation Study*

In this section, we perform ablation study on CUB dataset with VGGnet and GoogLeNet networks. We firstly evaluate the effect of each contribution over the baseline, then the effect of different perturbation values on localization and classification, and finally, how $\lambda_{CAM_{clean}}$ and $\lambda_{CAM_{adv}}$ influence the results of our method. The results are reportes in Tab. 6 and Tab. 7 and Fig. 7.

*Effect of AL and Entropy.* As shown in Tab. 6, with VGGnet backbone, using clean examples and adversarial examples improves greatly the baseline i.e 1.55% for classifcation and 14.92% for localization. When we further apply the entropy, we improve the localization accuracy by 0.09% with a little drop in classification of 0.02%. With GoogLeNet backbone, using adversarial learning improves the baseline with a margin of 4.12% for localization, however it drops the classification accuracy by 1.73%. When the entropy is further applied we improve both the classification and the localization by 0.04% and 0.56%, respectively

18

| Method | top1 cls-err | top1 loc-err |
|---|---|---|
| VGGnet-CAM [1] | 23.4 | 55.85 |
| VGG + Adversarial learning | **21.85** | 40.93 |
| VGG + Adversarial learning + Entropy | 21.87 | **40.84** |
| GoogLeNet-CAM [1] | 26.2 | 58.94 |
| GoogLeNet + Adversarial learning | 27.93 | 54.82 |
| GoogLeNet + Adversarial learning + Entropy | **27.89** | **54.26** |

Table 6: The effect of adding adversarial training and entropy loss to the baseline CAM.

*Adversarial Attacker Strength.* We show the effect of Projected Gradient Descent (PGD) [19] attackers with different perturbation values. we train both VGG network and GoogLeNet network on CUB dataset and as [24], we use perturbations $\epsilon$ ranging from 1 to 4 with an iteration of $n = \epsilon + 1$, except when $\epsilon = 1$ we set the number of iteration to 1. As shwon in Fig. 7, the bigger the perturbation, the higher is the error for both classification and localization. We get the best results with $\epsilon = 1$ and $n = 1$. This is obvious as our goal is not to build a robust model, but using adversarial learning as a way to activate more relevant features in the image.

*Regularization factors.* We further show the effect of the regularization factors $\lambda_{CAM_{clean}}$ and $\lambda_{CAM_{adv}}$ on the method, As shown in Tab. 7, with VGGnet backbone, we got the best results with $\lambda_{CAM_{clean}} = 1$ and $\lambda_{CAM_{adv}} = 0.01$. By selecting the right values for $\lambda_{CAM_{clean}}$ and $\lambda_{CAM_{adv}}$, entropy improves the localization accuracy.

## 5. Conclusion

In this paper, we proposed to take advantage of adversarial learning and entropy to improve WSOL performance. To do this, we train the model with clean examples and adversarial examples. By introducing some perturbations to the images, adversarial examples act as data augmentation and regularize the network, resulting in the activation of more relevant features. Furthermore,
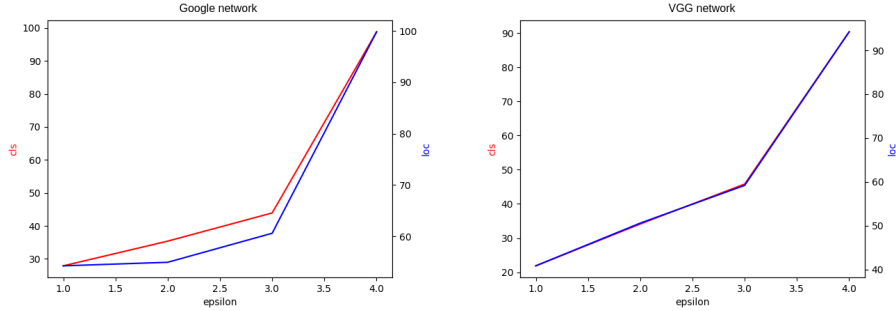
Figure 7: The effect of different perturbation strength on the classification and localization accuracy.

| Method | top1 cls-err | top1 loc-err |
|---|---|---|
| $\lambda_{CAM_{clean}} = 3,\ \lambda_{CAM_{adv}} = 1$ | 22.59 | 41.94 |
| $\lambda_{CAM_{clean}} = 1,\ \lambda_{CAM_{adv}} = 0.01$ | **21.87** | **40.84** |
| $\lambda_{CAM_{clean}} = 0.1,\ \lambda_{CAM_{adv}} = 0.01$ | 22.14 | 41.54 |
| $\lambda_{CAM_{clean}} = 0.01,\ \lambda_{CAM_{adv}} = 0.002$ | 21.88 | 41.23 |
| $\lambda_{CAM_{clean}} = 0.001,\ \lambda_{CAM_{adv}} = 0.0002$ | 22.11 | 41.09 |

Table 7: the effect of $\lambda_{CAM_{clean}}$ and $\lambda_{CAM_{adv}}$ on the results.

applying entropy minimization on the CAMs generated by the network, guides it during the training by forcing the pixels considered not relevant by the model to have a low entropy, and hence a higher prediction. Extensive experiments demonstrate that our EGA model obtained state of the arts on the three most used benchmark CUB, ILSVRC and OpenImages datasets.

## References

## References

[1] B. Zhou, A. Khosla, A. Lapedriza, A. Oliva, A. Torralba, Learning deep features for discriminative localization, in: The IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2016.

[2] X. Zhang, Y. Wei, J. Feng, Y. Yang, T. Huang, Adversarial complementary learning for weakly supervised object localization, in: IEEE CVPR, 2018.

[3] K. Kumar Singh, Y. Jae Lee, Hide-and-seek: Forcing a network to be meticulous for weakly-supervised object and action localization, in: The IEEE International Conference on Computer Vision (ICCV), 2017.

[4] J. Choe, H. Shim, Attention-based dropout layer for weakly supervised object localization, in: The IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2019.

[5] X. Zhang, Y. Wei, G. Kang, Y. Yang, T. Huang, Self-produced guidance for weakly-supervised object localization, in: The European Conference on Computer Vision (ECCV), 2018.

[6] H. Xue, C. Liu, F. Wan, J. Jiao, X. Ji, Q. Ye, Danet: Divergent activation for weakly supervised object localization, in: The IEEE International Conference on Computer Vision (ICCV), 2019.

[7] S. Yang, Y. Kim, Y. Kim, C. Kim, Combinational class activation maps for weakly supervised object localization, in: The IEEE Winter Conference on Applications of Computer Vision (WACV), 2020.

[8] S. Yun, D. Han, S. J. Oh, S. Chun, J. Choe, Y. Yoo, Cutmix: Regularization strategy to train strong classifiers with localizable features, in: The IEEE International Conference on Computer Vision (ICCV), 2019.

[9] F. Wan, P. Wei, J. Jiao, Z. Han, Q. Ye, Min-entropy latent model for weakly supervised object detection, in: The IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2018.

[10] B. Hakan, V. Andrea, Weakly supervised deep detection networks, in: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 2016.

[11] P. Tang, X. Wang, X. Bai, W. Liu, Multiple instance detection network with online instance classifier refinement, in: The IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2017.

[12] Y. Wei, Z. Shen, B. Cheng, H. Shi, J. Xiong, J. Feng, T. Huang, Ts2c: Tight box mining with surrounding segmentation context for weakly supervised object detection, in: The European Conference on Computer Vision (ECCV), 2018.

[13] F. Wan, C. Liu, W. Ke, X. Ji, J. Jiao, Q. Ye, C-mil: Continuation multiple instance learning for weakly supervised object detection, in: The IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2019.

[14] Y. Zhang, Y. Bai, M. Ding, Y. Li, B. Ghanem, W2f: A weakly-supervised to fully-supervised framework for object detection, in: The IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2018.

[15] J. Lee, E. Kim, S. Lee, J. Lee, S. Yoon, Ficklenet: Weakly and semi-supervised semantic image segmentation using stochastic inference, in: The IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2019.

[16] Z. Huang, X. Wang, J. Wang, W. Liu, J. Wang, Weakly-supervised semantic segmentation network with deep seeded region growing, in: The IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2018.

[17] X. Wang, S. You, X. Li, H. Ma, Weakly-supervised semantic segmentation by iteratively mining common object features, in: The IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2018.

[18] J. Choe, S. J. Oh, S. Lee, S. Chun, Z. Akata, H. Shim, Evaluating weakly supervised object localization methods right, in: The IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), 2020.

[19] A. Madry, A. Makelov, L. Schmidt, D. Tsipras, A. Vladu, Towards deep learning models resistant to adversarial attacks, in: International Confer-

ence on Learning Representations, 2018.

URL `https://openreview.net/forum?id=rJzIBfZAb`

[20] I. J. Goodfellow, J. Shlens, C. Szegedy, Explaining and harnessing adversarial examples, CoRR abs/1412.6572.

[21] T.-H. Vu, H. Jain, M. Bucher, M. Cord, P. Perez, Advent: Adversarial entropy minimization for domain adaptation in semantic segmentation, in: The IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2019.

[22] W. Wan, J. Chen, T. Li, Y. Huang, J. Tian, C. Yu, Y. Xue, Information entropy based feature pooling for convolutional neural networks, in: The IEEE International Conference on Computer Vision (ICCV), 2019.

[23] D. Tsipras, S. Santurkar, L. Engstrom, A. Turner, A. Madry, Robustness may be at odds with accuracy, in: International Conference on Learning Representations, 2019.

URL `https://openreview.net/forum?id=SyxAb30cY7`

[24] C. Xie, M. Tan, B. Gong, J. Wang, A. L. Yuille, Q. V. Le, Adversarial examples improve image recognition, in: The IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), 2020.

[25] T. Miyato, S.-i. Maeda, M. Koyama, S. Ishii, Virtual adversarial training: A regularization method for supervised and semi-supervised learning, IEEE Transactions on Pattern Analysis and Machine Intelligence PP. `doi:10.1109/TPAMI.2018.2858821`.

[26] S. Park, J. Park, S.-J. Shin, I.-C. Moon, Adversarial dropout for supervised and semi-supervised learning (2018).

URL `https://aaai.org/ocs/index.php/AAAI/AAAI18/paper/view/16322`

[27] A. Shafahi, M. Najibi, M. A. Ghiasi, Z. Xu, J. Dickerson, C. Studer, L. S. Davis, G. Taylor, T. Goldstein, Adversarial training for free!, in:

H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, R. Garnett (Eds.), Advances in Neural Information Processing Systems 32, Curran Associates, Inc., 2019, pp. 3358–3369. URL `http://papers.nips.cc/paper/8597-adversarial-training-for-free.pdf`

[28] S. Lee, D. Kim, N. Kim, S.-G. Jeong, Drop to adapt: Learning discriminative features for unsupervised domain adaptation, in: The IEEE International Conference on Computer Vision (ICCV), 2019.

[29] M. Everingham, S. M. A. Eslami, L. Van Gool, C. K. I. Williams, J. Winn, A. Zisserman, The pascal visual object classes challenge: A retrospective, International Journal of Computer Vision 111 (1) (2015) 98–136.

[30] K. Saito, D. Kim, S. Sclaroff, T. Darrell, K. Saenko, Semi-supervised domain adaptation via minimax entropy, in: The IEEE International Conference on Computer Vision (ICCV), 2019.

[31] Y. Grandvalet, Y. Bengio, Semi-supervised learning by entropy minimization, Vol. 17, 2004.

[32] S. Ioffe, C. Szegedy, Batch normalization: Accelerating deep network training by reducing internal covariate shift, ArXiv abs/1502.03167.

[33] C. Wah, S. Branson, P. Welinder, P. Perona, S. Belongie, The Caltech-UCSD Birds-200-2011 Dataset, Tech. Rep. CNS-TR-2011-001, California Institute of Technology (2011).

[34] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, L. Fei-Fei, ImageNet: A Large-Scale Hierarchical Image Database, in: CVPR09, 2009.

[35] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla, M. Bernstein, A. C. Berg, L. Fei-Fei, ImageNet Large Scale Visual Recognition Challenge, International Journal of Computer Vision (IJCV) 115 (3) (2015) 211–252. `doi:10.1007/s11263-015-0816-y`.

[36] T. Deselaers, B. Alexe, V. Ferrari, Weakly supervised localization and learning with generic knowledge, International Journal of Computer Vision 100. `doi:10.1007/s11263-012-0538-3`.

[37] K. Simonyan, A. Zisserman, Very deep convolutional networks for large-scale image recognition, in: International Conference on Learning Representations, 2015.

[38] D. Mahajan, R. Girshick, V. Ramanathan, K. He, M. Paluri, Y. Li, A. Bharambe, L. van der Maaten, Exploring the limits of weakly supervised pretraining, in: The European Conference on Computer Vision (ECCV), 2018.

[39] K. Simonyan, A. Vedaldi, A. Zisserman, Deep inside convolutional networks: Visualising image classification models and saliency maps, CoRR abs/1312.6034.