

Model	# of params	#of clusters	CUB-SS	CUB-PS	AwA-SS	AwA1-PS	FLO	Average
LDF(2018) [1]	426.4M	-	67.1	67.5	83.4	65.5	-	-
Resnet152	60.2M	-	66.9	67.3	81.0	67.5	64.0	69.3
Ours w/ SPN	61.0M	2	70.1	70.5	<b>83.7</b>	68.5	64.2	71.4
Ours	61.0M/42.5M	2	<b>70.5/66.5</b>	71.0/67.4	83.5/ <b>82.9</b>	<b>68.8/66.1</b>	<b>65.9/65.6</b>	<b>71.8/69.7</b>
	-	3	69.2/ <b>67.3</b>	<b>71.7/67.1</b>	82.4/82.6	66.3/66.5	65.8/64.7	71.1/69.4
	-	4	70.2/67.1	<b>71.3/67.6</b>	82.0/81.9	68.4/65.9	64.2/ <b>65.6</b>	71.2/69.6

Table 1: Zero-shot learning results on three benchmarks. The number of parameters is calculated for the CUB dataset. We report the results of our model **without/with** sharing the CNN parameters for the input image and the local patches.

Method	CUB			AwA1			AwA2			SUN		
	$A_{U \rightarrow T}$	$A_{S \rightarrow T}$	$H$	$A_{U \rightarrow T}$	$A_{S \rightarrow T}$	$H$	$A_{U \rightarrow T}$	$A_{S \rightarrow T}$	$H$	$A_{U \rightarrow T}$	$A_{S \rightarrow T}$	$H$
DEM [2]	19.6	57.9	29.2	32.8	84.7	47.3	30.5	86.4	45.1	20.5	34.3	25.6
RN* [3]	<b>38.1</b>	61.1	47.0	31.4	<b>91.3</b>	46.7	30.0	<b>93.4</b>	45.3	20.1	35.6	25.7
LDF* [1]	26.4	<b>81.6</b>	39.9	9.8	87.4	17.6	-	-	-	-	-	-
Ours*	36.7	71.3	<b>48.5</b>	<b>37.6</b>	87.1	<b>52.5</b>	<b>36.0</b>	84.3	<b>50.5</b>	<b>22.3</b>	<b>39.5</b>	<b>28.5</b>

Table 2: Generalized zero-shot learning results (%).  $H$  denotes the harmonic mean. \* means end-to-end training.

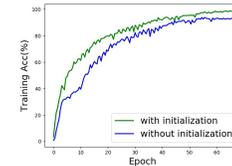


Figure 1: Training curve.

1 We first thank all reviewers for the valuable feedback.

## 2 Reviewer1

3 **Q1: Much more parameters.** To prove that the gain of performance is not totally from the higher capacity network,  
4 we conduct experiments using Resnet152 as the backbone with end-to-end finetune, which has a comparable amount of  
5 parameters to ours. As shown in Table 1, our model outperforms Resnet152 by 3.6%(71.8% v.s. 69.3%). Besides, our  
6 model has **significantly less parameters** than the best competing model LDF [1] while performing much better. We  
7 also can reduce the number of parameters by using the same CNN for the image and the part patches as you suggested,  
8 but the ZSL performance is slightly degraded roughly from 71% to 69% as shown in Table 1 (separated with slashes).

9 **Q2: The importance of weights initialization.** We present the training curves of our model with/without weights  
10 initialization in Fig 1. We see that initializing the attention layers speeds up the learning and finally achieves a greater  
11 accuracy. We will add more detailed analysis in the final version of the paper.

12 **Q3: The number of clusters.** As shown in Table 1, we increase the number of clusters to 4 and find little performance  
13 improvement. Besides, we observe more maps introduce the attention redundancy, i.e. maps attend to the same region.

14 **Q4: Results for generalized ZSL.** As shown in Table 2, our model outperforms the other SOTA models (based on  $H$ ).  
15 **Other comments.** The CNN pretrained to provide pseudo labels for clustering is the same backbone used in our model,  
16 otherwise it would give erroneous peak as you agreed. We will cite the relevant papers you suggest.

## 17 Reviewer2

18 **Q1: About Cropping network.** In fact, to obtain better representation for finer localized cropped region  $x_i^{part}$ , our  
19 method also utilizes the bilinear sampling to adaptively zoom the cropped region  $x_i^{part}$  to the same size with the original  
20 image. Concretely, for a point  $(i, j)$  of the zoomed region, its value  $x_{(i,j)}^{zoom}$  can be computed bilinearly combining the  
21 values of nearest four points in the cropped region. Formally,  $x_{(i,j)}^{zoom} = \sum_{\alpha, \beta} |1 - \alpha - \{i/\lambda\}| |1 - \beta - \{j/\lambda\}| x_{(m,n)}^{part}$ ,  
22 where  $m = [i/\lambda] + \alpha + z_x - z_s$ ,  $n = [j/\lambda] + \beta + z_y - z_s$ ,  $\alpha = 0$  or  $1$ ,  $\beta = 0$  or  $1$ ,  $\lambda$  is the upsampling factor,  $\lambda = t/t_s$  ( $t$   
23 is the size of the original image) and  $[\cdot]$  and  $\{\cdot\}$  is the integral and fractional part, respectively. We will add the detailed  
24 description in the final version of the paper. Spatial Transformer Network (SPN) is an alternative of our cropping net.  
25 When replacing it with SPN, we find the performance changes little as shown in Table 1.

26 **Q2: About triplet loss.** We agree that the normalization will change the relative distance of two points. There is a typo  
27 leading to misunderstanding in the paper. We actually use the normalized version of  $\phi$  in the embedding softmax loss  
28 so that only normalized features are considered and used in training and inference phases. Please refer to [4]. We will  
29 add more discussion in the final version of paper. **Other comments:** The local patches are going through the same  
30 CNN while the input image is going through a different CNN. We will mark it in the figure.

## 31 Reviewer3

32 The competing method LDF [1] used a single attention scheme and we have shown our superiority to it in both the  
33 model design and the performance. We will add more explanation about how the extracted features are used and add a  
34 reference to the appendix. Please kindly refer to our response to other reviewers.

## 35 References

- 36 [1] Li et al. Discriminative learning of latent features for zero-shot recognition. In *CVPR*, pages 7463–7471, 2018.
- 37 [2] Zhang et al. Learning a deep embedding model for zero-shot learning. In *CVPR*, pages 2021–2030, 2017.
- 38 [3] Sung et al. Learning to compare: Relation network for few-shot learning. In *CVPR*, pages 1199–1208, 2018.
- 39 [4] Wang et al. Normface: l2 hypersphere embedding for face verification. In *ACMMM*, pages 1041–1049. ACM, 2017.