

1 We would like to thank reviewers for recognizing our original contribution (R7), convincing experiments (R6, 7, 8) and  
2 clarity/reproducibility of the paper (R6, 8). We appreciate suggestions from R6, 7, 8 and will include these in the paper.

### 3 Response to Reviewer 5

4 – **C1:** “*lack of novelty: pretty similar to DAZLE*” **Answer:** We respectfully disagree. As mentioned by R7, our  
5 contribution is to generate attribute-based features and compose them for recognition of unseen and seen classes, which  
6 is original and has not been done before, including DAZLE.

7 – **C2:** “*how to classify unseen classes.*” **Answer:** We use the discriminative model  $p(y|\mathbf{H}, \mathbf{z})$  to compute probabilities  
8 of unseen classes given their semantic vectors  $\mathbf{z}$  (line 187) and classify a sample as its most probable unseen class.

9 – **C3:** “*... compared with more state-of-the-art methods.*” **Answer:** As mentioned by R6, 7, 8, our experiments are  
10 extensive. We have included most competitive methods with comparable settings to ours at the submission time.

11 – **C4:** “*What if the attributes of one unseen class are not shared with seen classes?*” **Answer:** Please notice that ALL  
12 zero-shot learning works rely on sharing attributes between seen and unseen classes. Without sharing common attributes  
13 between seen and unseen classes, there is no way to transfer knowledge to unseen classes for recognition.

### 14 Response to Reviewer 6

15 – **C1:** “*The improvements ... on CUB and DeepFashion are limited.*” **Answer:** Please notice that we significantly  
16 improve zero-shot accuracy by at least 3.2% and 2.2% on DeepFashion and CUB, respectively, compared to other  
17 methods (lines 252-253) while having no extra learnable parameter w.r.t. DAZLE.

18 – **C2:** “*performance of ... seen categories is not so good.*” **Answer:** Please notice that high performance on seen classes  
19 is not the main goal of zero-shot learning as it can cause seen class bias (low unseen class performance). Our method  
20 achieves competitive seen class accuracies while obtaining the best unseen accuracies, hence, high harmonic means.

### 21 Response to Reviewer 7

22 – **C1:** “*Providing e.g. pseudo code describing the whole  
23 process ... training process*” **Answer:** Thanks for the  
24 suggestion. We will include the shown Algorithm 1.

25 – **C2:** “*miss ... work relying on class similarity graphs*”  
26 **Answer:** Thanks. We will discuss these works, if ac-  
27 cepted. Compared to the reported numbers from the most  
28 comparable work (Ding et al., CVPR 2019), our method  
29 outperforms this work in harmonic mean by a significant  
30 margin of 25% on CUB, AWA2 and 5% on SUN.

31 – **C3:** “*Is a feature vector constructed per attribute? ... Is  
32 S sampled such that all attributes are present?*” **Answer:**  
33 As mentioned in Remark 1, the composed dense features  
34 consist of a feature vector per attribute which is required

35 for  $p(y|\mathbf{H}, \mathbf{z})$  in DAZLE. Although we try to include samples from various classes having different attributes (line 153,  
36 242), we cannot guarantee  $S$  contains all present attributes without extra annotations of present attributes in images.  
37 However, our framework can compose features from any set  $S$  by solving Eq (10) even with missing attributes in  $S$ .

38 – **C4:** “*How ... no Comp relate to the DAZLE model?*” **Answer:** Please notice that they are different. The No\_Comp  
39 variant is trained via cross-entropy loss while DAZLE [10] is trained with an extra self-calibration loss (lines 67-68).

40 – **C5:** “*Could the approach been extended to ... learned attributes?*” **Answer:** For learned attributes, we can factorize  
41 attribute representations into common components (via PCA) which can be used as the building blocks for composition.

### 42 Response to Reviewer 8

43 – **C1:** “*semantic vectors  $\mathbf{z}^c$  are available at training time for both seen \*and\* unseen classes ... [15, 16, 17], ... goes  
44 against the philosophy of zero-shot learning.*” **Answer:** To be comparable with [15 - 17], we follow their setting which  
45 enables our model to compose unseen class features during training, thus the model can learn the testing distributions  
46 with both seen and unseen classes (lines 70-72). Please notice that without the availability of unseen class semantics at  
47 training time, we cannot use self-composition to alternate between training a classifier and composing features.

48 – **C2:** “*how they obtained the results of Table 1 and Table 2 (left) for prior works?*” **Answer:** Thanks. We will include  
49 the following clarification: On DeepFashion, we run each baseline using their released codes with their default settings.  
50 On the remaining datasets, we use the performances reported in their papers to ensure their best performances.

51 – **C3:** “*Do all these works use a ResNet 101 backbone*” **Answer:** Except for SMA using VGG19 and LFGAA combining  
52 VGG19, GoogleNet, and ResNet101, all remaining baselines use ResNet101. We will clarify this in the paper.

53 – **C4:** “*does one really need to restrict the set of possible features to semantically-related samples?*” **Answer:** Please  
54 notice that the related sample set  $Q(\mathbf{z})$  is required for the sampling probability  $p(\mathbf{H}|\mathbf{z})$  and is computed via nonnegative  
55 OMP (lines 171-173 and Algorithm 1). Without related samples, we cannot solve Eq (10) efficiently through sampling.

56 – **C5:** “*Lines 274-275 ... Can you please clarify?*” **Answer:** Please notice that we set the attribute values of each sample  
57  $\mathbf{z}^i$  to its class semantic  $\mathbf{z}^{y_i}$  (lines 239-240). However, due to occlusion, a sample has many missing attributes compared  
58 to its  $\mathbf{z}^{y_i}$ , thus relying only on  $\mathbf{z}^{y_i}$  without  $\mathbf{H}_i$  would compose features lacking many discriminative attributes.

59 We hope our responses answer the questions and kindly ask the reviewers to raise their scores in light of our responses.

---

#### Algorithm 1 Self-Composition

```
1: Input: Training set  $D$ , DAZLE model, ResNet101 backbone
2: Initialize discriminative model  $p(y|\mathbf{H}, \mathbf{z})$  from DAZLE
3: for  $t = 1, \dots, N_{iteration}$  do
4:   Sample a set  $S \subset D$  with an equal number of samples per
   class (lines 153 and 242)
5:   Extract region features  $\{\mathbf{f}_i^r\}_{r=0}^R$  for  $i \in S$  via ResNet101
6:   Extract dense features  $\mathbf{H}_i$  for  $i \in S$  via DAZLE model
7:   Construct  $Q(\mathbf{z}^u)$  for  $u \in C_u$  via solving Eq (8) with OMP
8:   Compose features for  $u \in C_u$  based on  $Q(\mathbf{z}^u)$  via Eq (10)
9:   Update  $p(y|\mathbf{H}, \mathbf{z})$  on real/composed features via Eq (11)
10: end for
11: Output: Optimal discriminative model  $p(y|\mathbf{H}, \mathbf{z})$ 
```

---