

1 We thank all the reviewers for their valuable comments and positive feedback: R1) Indicating an important aspect of
2 the zero-shot task; R3) Rigorous experiments; R4) A novel method. All the reviewers are satisfied with our paper’s
3 contributions, good performance, and clear presentation. Codes for reproducing all the experiments will be released.

4 **Response to Reviewer-1.**

5 *R1.1 Comparisons with other papers.* We note that previous works [2,27,48] experiment with different settings and
6 datasets, which makes it hard to conduct unified and direct comparisons. We also considered ZS3Net[2] for direct
7 comparison, yet found that its setting may be different from ours. We assume no knowledge of unseen classes in
8 training, while ZS3Net learns implicitly with supervision of unseen classes. We would like to refer reviewers to Issue
9 #1, #4, #6 raised (by other researchers) in the official Github repo of ZS3Net for contrasting the settings. To ensure
10 fair comparisons in our work, for the GMMN we adopted the implementation from ZS3Net. And all the methods are
11 trained similarly with independently selected hyper-parameters for each.

12 *R1.2 How our method works.* As motivated in Fig.1, in zero-shot learning, optimizing over the noisy/abnormal samples
13 of seen classes may result in models with a biased visual-semantic mapping, thus the inferred classifier for unseen classes
14 may be less reliable. In the scenarios of semantic segmentation, the noise appears in two levels: one is at image-level,
15 e.g. some van cars look like ‘bus’ but are labelled as ‘car’; the other is at pixel-level, e.g. pixels near boundaries are
16 hard to distinguish and easily mis-annotated. Our uncertainty-aware model is proposed to address such challenges.
17 We observed that the image-level learning path is critical to the accuracy of segmentation; yet, at the same time, the
18 pixel-level learning path plays as a complementary role in learning refined segmentation details.

19 *R1.3 How the uncertainty-aware learning works.* Due to the noisy data collection and annotation process in practice,
20 data samples may have different levels of uncertainty. For example, an abnormal sample is expected to make a less
21 confident prediction than a typical sample, and pixels in the object centers are expected to be more confident than
22 boundary pixels. The level of noise introduced for individual samples is called Heteroscedastic uncertainty and can be
23 formulated as a random variable and estimated nonparametrically along with the learning of deep neural networks [12,
24 19, 28, 32, 38]. Specifically, in deep networks for classification or regression, the model is trained to learn a feature
25 mapping such that the features are distributed in the form of a certain exponential family, with a unique variance and
26 mean for all samples. Optimizing the feature mapping over noisy/abnormal samples may result in a biased distribution.
27 Yet, explicitly modeling the uncertainty σ for each sample, as in our formulation, allows the model to adaptively adjust
28 the variance for individual samples. As a result, all the data samples can be properly accounted for by the model with
29 less bias toward the abnormal/noisy samples. During the optimization process, the abnormal/noisy samples are typically
30 mapped far from the majority of the samples, thus leading to high uncertainty. We will clarify this in the paper.

31 *R1.4 Parameter selection.* We select hyper-parameters based on a validation set split from the seen-class training set.

32 *R1.5 Eq 7: for sake of...* We also tried with L-2 loss, but found the accuracy drops in all settings.

33 *R1.6 Table 1: DeViSE results...* One main reason may be that DeViSE maps the visual feature into the low-dimensional
34 semantic space for nearest neighbor search. Such a process may shrink useful information and aggravates the Hubness
35 problem in zero-shot learning. Moreover, in this work, we test with a much larger unseen class set than in previous
36 works, thus further increasing the difficulty of classification for DeViSE.

37 *R1.7 Table 2: Adding Random...* We removed these two lines from this table to align the height of the figure in the right.
38 On ADE20k data, the mIoU (in %) of the unseen classes for Blank/Random is 11.2/ 5.8 (K=25), 9.3/ 5.1 (K=50), and
39 9.2/ 4.7 (K=75), which is much lower than our final accuracy. We will add this back.

40 *R1.8 Fig 3b: what is the intuition...* We guess that R1 refers to Fig 3c. We found this is a dataset-specific phenomenon
41 and observed consistent improvements of unseen-class performance for all the methods on PC30 under this setting.
42 The reason may be that removing the unseen-class images results in smaller but more compact training set for seen
43 classes, thus leading to better visual-semantic mappings for inferring model for unseen classes. We also observed that
44 the overall mIoU for our final model drops from 36.5% to 28.6% when removing all the unseen-class images for PC30.

45 *R1.9 Other comments.* We appreciate the detailed and very helpful comments. Due to the page-limit, we can only
46 address part of them in this rebuttal. We will carefully address all the comments in a revised version of our paper.

47 **Response to Reviewer-3.**

48 *R3.1 About pixel and observation noise.* Yes, noise exists at both the image- and pixel-level. Please see our response to
49 R1.2 and R1.3 for details.

50 *R3.2 Baseline classifying pixels independently.* Thanks for the suggestions. With only the pixel-level path, the model
51 achieves mIoU for unseen classes with 13.5% for PC30, 8.8% for PC156, and 7.9% for ADE75. Compared to the
52 results in Tab.3, we found that simply adding the image-level path can improve about 0.4%, 1.9%, and 5.6% for PC30,
53 PC156, and ADE75 respectively. We will update our paper.

54 *R3.3 Generator-based techniques.* We denote [2] as GMMN in our experiments. Please also see our response to R1.1.

55 **Response to Reviewer-4.**

56 *R4.1 Results for different uncertainty/noise level.* Thanks for the suggestions. To evaluate the effects of different levels
57 of noise, we randomly shuffle the class labels for a training sample. The unseen-class mIoU (in %) for baseline/U-loss
58 under different probability on PC-156 are 11.7/13.3 (p=0.0), 11.1/12.8 (p=0.1), 10.9/12.9 (p=0.25), 9.6/12.1 (p=0.5),
59 which shows that uncertainty-aware learning achieves better robustness. We will update the paper.