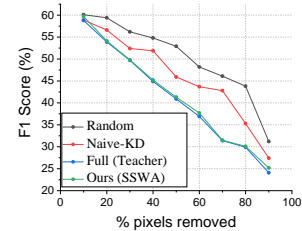1    We thank the reviewers for their constructive feedback. We are encouraged that all reviewers found our observation
2  and motivation - attribution map distortion while compressing networks - to be novel and interesting [R1,R2,R3,R4]
3  and our approach to be intuitive [R1], readily applied [R2], and evaluated with convincing experiments [R1,R3]. We
4  address the reviewer comments as much as space permits and will include all feedback in the final version.

5  **[R1,R2] Evaluation of predictive performance and at-**
6  **tribution score with *Error Bound*.** We reported the (multi-
7  label) predictive performance measure in main the text with
8  *mean-average-precision (mAP)*. We apologize for not elab-
9  orating this [R1,R2]. Also, "accuracy" in Appendix table
10 3,4 refers to the ImageNet top-1 accuracy [R1]. Moreover,

| | Predictive Performance | | Attribution Score | |
|---|---|---|---|---|
| Method | mAP | F1 Score | AUC | Point Acc |
| Full (Teacher) | 91.79±0.16 | 78.44±0.23 | 88.68±0.17 | 80.16±0.16 |
| Naive-KD | 81.44±0.19 | 62.71±0.28 | 80.51±0.25 | 69.18±0.17 |
| EWA | 82.31±0.22 | 63.16±0.31 | 84.23±0.22 | 79.19±0.22 |
| SWA | 83.63±0.19 | 64.94±0.24 | 87.73±0.19 | 79.91±0.21 |
| SSWA | **84.37±0.26** | **66.27±0.35** | **88.13±0.22** | **80.14±0.24** |

11 we will add one more predictive performance measure: F1 score. Following the [R2]'s comments, we additionally
12 report the performances with standard variation through the 5 same experiments to be more convincing results. In above
13 the table, we report the predictive performances and attribution scores of various methods including our methods for
14 knowledge distillation with vgg/4. For other compression cases, we will incorporate in the final version.

15 **[R1,R2,R3,R4] Other evaluation of attribution map & Reliability of Grad-Cam.** Fol-
16 lowing R1's suggestion, we will add the perturbation metric, RemOve And Retrain
17 (ROAR) [A], to evaluate the how well attribution maps from compressed networks explain
18 the model behavior. This test is as follows: remove the top-$k$ pixels of an image ranked
19 by the attribution map for the entire training data, and retrain a classifier on this perturbed
20 dataset. If the attribution map accurately represents the importance of the pixels, the
21 classifier must have lower predictive performance. Here, we show a quick experiment



22 using vgg-11 classifier. **(a)** In the right graph, all Grad-Cam perturbations (from different
23 models) were able to lower the F1 score more than random perturbations, which verifies that Grad-Cam indeed reflects
24 a model's decision making process. **(b)** The student trained with our method scored almost in par with the full network.
25 This indicates that the attributions (which reflect a model's decision process) are indeed preserved by our method.

26 **[R4] Preserving differentiable attribution map also preserves**
27 **other attributions.** We observed the deformation of various attribu-
28 tion maps in Appendix A. Following R4's suggestion, we evaluated

| AUC/PointAcc | Excitation Bp | LRP$_{\alpha=1\beta=0}$ | RAP |
|---|---|---|---|
| Full (Teacher) | 84.2/74.8 | 85.3/65.5 | 84.5/69.5 |
| Naive-KD | 76.3/66.3 | 79.6/53.4 | 80.9/56.9 |
| SSWA (Ours) | **82.3/71.2** | **82.5/64.1** | **83.5/65.7** |

29 other attribution maps for the model using our method. We observe that it also helps preserving other attribution maps.
30 Since most attribution maps [14, 15] are generated with gradients and activations, preserving one may help others.

31 **[R2,R3,R4] Technical novelty.** In [4], the authors introduce attention map transfer and gradient transfer. Although
32 our method and theirs share certain aspects, we believe that there are enough differences to them. **(a)** The problem of
33 focus is distinct. [4] focuses on boosting the *predictive performance* of a student network, while our method focuses
34 on preserving the *attribution* power of a network while being compressed. **(b)** The function of the loss functions
35 are different. Attribution map used in attention matching in [4] lacks label-specific attribution information since all
36 activation maps are equally weighted and aggregated. In other words, this form of attribution map may tell where the
37 network is looking, but it holds no "meaning". Thus, this regularizer may teach the student how to look and distinguish
38 objects, but does not pass on the information of "what" and "why" it should look at a certain region. Gradient matching
39 of [4] does hold some label-specific information. However, since they only match the gradients at the input level, the
40 information in the intermediate layers (and thus their decision processes) are not appropriately transferred while our
41 method is able to transfer information from intermediate layers by collapsing the channels.

42 **[R2,R3] Motivation of stochastic matching.** The intent behind our stochastic matching regularizer is to facilitate
43 the transfer of relevant information and prevent overfitting. Several recent works utilize this concept to boost the
44 performance of knowledge transfer between a teacher and a student. In [23], they encourage the teacher-student
45 information transfer by using a gaussian dropout to maximize the mutual information between teacher and student
46 features. Injection of stochasticity can also be found on other similar fields such as continual learning [B] and domain
47 adaptation [C]. Based on the empirical evidences of performance increase presented by the literature mentioned above,
48 we believe that it is plausible to motivate our stochastic matching method.

49 **[R3] Hypothetical explanation for the attribution deformation problem.** As a network is compressed, it must cram
50 its decision procedures(information) inside a smaller memory. If so, it would be unable to use "standard" decision
51 procedures, but must resort to using shortcuts and inklings that means less to humans. Thus, its decision procedures
52 would become harder to interpret which is reflected in its distorted attribution map.

53 **Clarity and minor details.** [R2] Is the "full net" finetuned or trained from scratch?: *Fine-tuned.* [R3]··· is Grad-Cam
54 used for some other purpose?: *the qualities of Grad-CAM maps extracted from EWA, SWA, SSWA were evaluated.* [R3]
55 We apologize for the mis-reporting: In Appendix Table3, the 'Full Network''s AUC should be changed: *75.92 → 81.64*

[A] Hooker, Sara, et al. "A benchmark for interpretability methods in deep neural networks." NeurIPS. 2019.; [B] Lee, et al. "Overcoming Catastrophic Forgetting by Incremental Moment Matching", NeurIPS. 2017.; [C] Saito, et al. Adversarial Dropout Regularization, ICLR. 2018.;