

1 We would like to thank all the reviewers for their time and for giving our paper this encouraging feedback.

2 **Reviewer #1** mentioned the issue of transferability of the student features when trained without any real data, which
3 is also something we consider interesting. Since submission we tested our student on the ImageNet component from
4 the CINIC-10 dataset (<https://arxiv.org/abs/1810.03505>), corresponding to a subset of ImageNet images chosen and
5 resized to match CIFAR. We observe that the student generalizes equally well to its teacher, i.e. that the difference in
6 performance between teacher and student remains the same under such dataset shift. This is not as much of a shift as
7 transferring to a segmentation task (in general classification features generalize poorly to segmentation tasks even in
8 normal circumstances). But in general we find that the student’s representation shows similar transferability to that of
9 the teacher, because of the student’s tendency to directly match the decision boundaries of the teacher. However we
10 cannot make any claim of increased transferability at this stage.

11 Vaccination is also an interesting thought. In practice the nature of the problem forces the student to match the teacher
12 on the entire input space regardless of how it was trained, but having a teacher agnostic to out of distribution data could
13 reduce the task-relevant information that leaks out to the student. On the other hand, depending on the vaccination
14 approach it could be that the agnosticism of the teacher helps identifying samples close to the training data more easily,
15 because they would have a specific entropy signature (high confidence). We have not pushed further on vaccination at
16 this stage because it seems to be a use-case that is harder to justify in practice: someone releasing a trained model for
17 people to use typically wouldn’t take measures against people compressing it.

18 We agree with the reviewers that additional large scale datasets would be complementary to this work. However we also
19 note that many of the competing papers we cited are limited to MNIST experiments due to the difficulty of the problem,
20 so our tests are more extensive than many. In practice, a limitation to our current model is computational cost, since
21 our distillation trains two networks instead of one, and several gradient steps per batch are taken for the student. This
22 amounts to a handful times more computational cost than training a vanilla network, and our limited computational
23 resources have made ImageNet scale experiments a challenge. This is a focus of future work; however those working in
24 this field need to know of methods that are effective on the smaller problems so we can focus our efforts on scale in the
25 right direction. Hence we think it is important to not wait for huge-scale experiments before releasing this work.

26 In particular in future work we plan to investigate more closely the relationship between the size of the teacher’s training
27 set and the difficulty of zero-shot distillation. Extra challenges relating to large scale datasets may include the larger
28 number of classes, some of which may not be sampled enough by the generator.

29 We thank **Reviewer #2** for his comments on layout, and have modified our paper in light of that for the camera ready
30 version.

31 **Reviewer #4** correctly and pertinently noticed that in Figure 4 we took adversarial steps on the students rather than the
32 teacher. Attacking the students rather than the teacher allows us to show 4 curves per dataset instead of 3, and as such it
33 conveys a bit more information. For instance it shows that it is easier to cross decision boundaries on our zero-shot
34 student than on the baseline student when stepping on each network respectively (orange vs red curve per row). Another
35 possible motivation comes from considering the case that the teacher and student are sets of linear decision boundaries.
36 If we assume that the students is a smaller set (less capacity) mostly contained within the teacher set, then crossing
37 decision boundaries of the student should correspond to crossing decision boundaries of the teacher, to a greater extent
38 than the other way around.

39 Lastly we note that the batchnorm layers of the WRNs in our paper were initially not trained (gamma and beta fixed).
40 This was a mistake because a trainable batchnorm is the most common practice for WRNs in the literature, and it
41 facilitates the training procedure. We re-ran all the experiments in the paper and obtained a 1 to 2% boost to most of
42 our accuracies, and reduced the Transition Error (TE) on CIFAR-10 by a third. For the sake of reproducibility we will
43 ensure the camera ready version has these updated values, so that all values match our latest published code.

44 Again we thank the reviewers for their time and hope that we have answered their questions.