

1 We thank the reviewers for thoughtful feedback. In this paper we introduce Augerino, a method that can automatically  
2 attain the strength of group equivariant networks through learning data augmentations. Augerino is a flexible method that  
3 can be combined with any pre-existing architecture. We note that Augerino does not require one to specify invariances  
4 a priori, merely a group transformation. We show Augerino is capable of **learning** augmentations, from training data  
5 alone, which boosts accuracy when applied at test time over a wide range of tasks. We want to emphasize that Augerino  
6 is a distinctive contribution: other works on data augmentation are not geared towards learning equivariances at training  
7 time and would not even directly apply to most of the applications considered in our paper.

8 **R1:** Thank you for the constructive review, we appreciate the feedback and will incorporate the suggestions accordingly.  
9 With respect to the discussion of the mathematics, we commit to making the mathematical statements in Sections 3.1  
10 and 3.2 more precise, we will switch from the notation of the transformation and distribution of  $T_\phi, Q_\theta$  to  $g, \mu$ . Our use  
11 of uniform distribution should be made more precise, and that it could depend on the parameterization if not careful. We  
12 will make it clear that the distribution is meant to be the Haar measure on the group, and that the appropriate condition is  
13 left translation invariance  $d\mu(hg) = d\mu(g)$ . The transformation  $T_\phi$  needs to act linearly so that we can pull it out of the  
14 expectation in Eqs. 9 – 10. Unlike with learning invariances, the extension to equivariance in Section 3.1 does require  
15 the group structure and is not applicable to more general sets of transformations which may lack inverses or closure, and  
16 we will add this point of discussion in the text. While this requirement is an idealization and limitation in some ways  
17 (cropping, boundary effects), we feel that this simplification is justified. Note that the invariant model only requires an  
18 invariant measure, but not necessarily the existence of inverses, and thus includes transformation semi-groups. We will  
19 mention the deep scale-spaces paper in this context. Additionally, thank you for the pointer to the Barnard and Casasent  
20 paper, this is a great early example of the utility of invariance.

21 **R2:** While we appreciate the pointers to test-time augmentation, which we are happy to reference, we want to be clear  
22 that these papers are complementary to Augerino and do not diminish our contribution. Indeed, most of these works  
23 would not even directly apply to the applications we consider. Augerino is distinctive in that it **learns** an appropriate  
24 distribution over augmentations at **train-time**. We motivate the utility of Augerino in Section 4, demonstrating both  
25 why one may want to learn a distribution over augmentations, rather than just applying them.

26 We also note that for Augerino to be applicable we need only the invariance to be represented by a group transformation.  
27 In Section 4 we do know the correct invariances, but these experiments serve to show that when we do know the  
28 correct distribution over invariances we can recover the ground truth, and are in no way meant to show the limiting  
29 case performance of the method. In Sections 5 and 6 we do not have a known target set of invariances yet Augerino  
30 outperforms competing methods. While in some cases the set of invariances will never be known, if we do believe  
31 that invariance is a trait we wish our model to have then group transformations are a promising path forward as they  
32 encompass a broad set of transformations that has been shown to improve performance on many well studied tasks; for  
33 examples see [6, 7, 8, 12].

34 **R3:** We thank the reviewer for the constructive comments and ideas for extensions. The model used in Section 5 is an 8  
35 layer deep convolutional neural network with a max channel width of 256. We will be sure to include the details of all  
36 architectures in the appendix of the camera ready.

37 We have not observed any evidence that Augerino is susceptible to loss of performance in using alternate network  
38 architectures. One of the key strengths of Augerino is its ability to be combined freely with any model. Throughout the  
39 experiments we use a range of different convolutional and feed forward neural networks, chosen only to be appropriate  
40 for the given task not for their compatibility with Augerino. In each case (the demonstrations of Section 4, and the  
41 benchmark improvements of Sections 5 – 8) Augerino performs as expected and improves accuracy over baselines.

42 Finally, we will incorporate larger datasets in the camera-ready submission. Preliminary results suggest Augerino will  
43 perform well on larger datasets such as CIFAR100, however we avoid presenting any incomplete results here.

44 **R4:** Thank you for your thoughtful comments. *Assumption of hyperparameters:* The parameters  $a$  and  $b$  are *learned*  
45 *parameters of the model* and not hyperparameters; the reparameterization trick gives us gradients with respect to these  
46 parameters. Therefore  $a$  and  $b$  can learn to accommodate an arbitrary range, and we do not find performance sensitive to  
47 their initialization. We will clarify. *Learning transformations:* as discussed in response to reviewer 2, the experiments  
48 in Section 4 highlight our ability to uncover the ground-truth augmentations. In later sections we show Augerino can  
49 outperform competing methods in which desired transformations are unknown.

50 *Remaining points:* The augmentation in Eq. 12 is done through computing the exponential map as the solution to a  
51 differential equation which allows back-propagation to flow through the affine transformations. In Section 5 the fixed  
52 augmentation indicates a standard set of transformations applied to the training data, we do expect fixed augmentation  
53 to perform worse, that is the central motivation to this work. We will be sure to include these details and clarify all  
54 notation in Section 3 in the camera-ready.