

1 Author response for “Fixing the train-test resolution discrepancy”

2 We thank the reviewers for their constructive feedback on the paper. We will take into account their comments on
3 presentation and typos. Here we answer their main questions and comments.

4 **R1: The results presented by the paper do not seem quite significant because just one dataset (ImageNet-2012)
5 and two models (ResNet and PNASNet) are used. In addition, are the results shown significant?** We will add
6 experiments on additional models and datasets.

7 In particular, we have evaluated our approach for transfer learning for low-resource and/or fine-grained classification.
8 Following common practice, (1) we initialize the model with weights obtained on ImageNet and (2) fine-tune it on
9 the new dataset, which provides the baseline performance. Then (3) we use our method, i.e. we fine-tune the last
10 batch norm and the fully connected layer with a higher resolution. Table 1 shows that for all transfer-learning tasks
11 our method obtains a gain in accuracy of +0.4% to +1.3% absolute. All those results are obtained on competitive
12 benchmarks where tens to hundreds of researchers have evaluated their methods, gains of around 1% accuracy are
13 therefore significant.

14 Finally, we applied our method to a very large ResNeXt-101 32x48d from [Mahajan et al. ECCV’18], available online.
15 We improved their top-1 Imagenet accuracy from 85.4% to 86.4%, which is the new state of the art on ImageNet.

16 **R2: why not just fix the training such that there is no discrepancy, as opposed to changing the size for test and
17 finetuning?** We have tried other approaches to fix the discrepancy, which are not discussed for space reasons. The
18 approach we presented consistently achieved the best performance while being simple and reducing the training time.

19 Regarding R2’s proposal, since the discrepancy is mainly due to the random resized crop data augmentation, we replaced
20 it with a fixed size random crop and resize to $K_{\text{train}}=224$ pixel. When testing at $K_{\text{test}}=224$ we obtain 74.6% accuracy.
21 The best test resolution is $K_{\text{test}}=256$: 75.2% accuracy, so there is an improvement at a slightly higher resolution.
22 However, the result is significantly below the baseline accuracy obtained with random resize crop and without any
23 resolution change (76.2%). This shows that it is important to keep a distinct data augmentation between train and test.

24 **R2: Line 110-111 derives $f = \sqrt{HW}$, which does not seem to be right since k doesn’t include the sensor size.**
25 As often done for compactness (see [Hartley & Zisserman] eq. 5.9 or 6.9 depending on the edition), we expressed
26 the focal length in unit of pixels; namely, the camera projection equation is $x = W_{\text{pixels}}^{\text{sensor}} / W_{\text{mm}}^{\text{sensor}} f_{\text{mm}} X/Z$ which
27 simplifies to $x = f X/Z$ by setting $f = f_{\text{pixels}} = W_{\text{pixels}}^{\text{sensor}} / W_{\text{mm}}^{\text{sensor}} f_{\text{mm}}$. We will clarify this in the final version.

28 **R3: what is a good ratio between the data augmentation hyperparameter(s) and the train/test size ratio? In
29 other words, a general rule of thumb for practitioners in image recognition on the ImageNet dataset.** Thank you,
30 we will add a “best practices” paragraph. It is hard to define a single ratio that works for all resolutions because the
31 difference in statistics between the activations in the neural network does not vary linearly, eg. because of the padding
32 of convolutional layers. An important effect to consider is the interaction between the input resolution and the size
33 of the feature map before spatial pooling. There are key resolutions (largest size for a given feature map size) that
34 provides local performance minima. For example, for a resnet-50 this occurs each time the resolution is a multiple of
35 32. Hence, a practical approach is to try out a few resolution steps above the training resolution, ie. {256, 288, ..., 384}.
36 This approach applies both with and without fine-tuning, but fine-tuning improves the effect.

37 **R3: And, are the observations generalizable? If I transfer the ImageNet model to a specific task (e.g., CUB),
38 will the findings in this paper be useful? If the answer is yes, how can it be useful?** Yes, please see our answer
39 to R1.

Dataset	Model	Baseline	Ours	State-Of-The-Art Models		
iNaturalist 2017	SENet-154	74.1	75.4	IncResNet-V2-SE	[Horn <i>et al.</i> ArXiv’17]	67.3
Stanford Cars	SENet-154	94.0	94.4	EfficientNet-B7	[Tan & Le, ArXiv’19]	94.7
CUB-200-2011	SENet-154	88.4	88.7	MPN-COV	[Peihua Li <i>et al.</i> ArXiV’19]	88.7
Oxford 102 Flowers	IncResNet-V2	95.0	95.7	EfficientNet-B7	[Tan & Le, ArXiv’19]	98.8
Oxford-IIIT Pets	SENet-154	94.6	94.8	AmoebaNet-B (6,512)	[Yanping Huang <i>et al.</i> ArXiv’18]	95.9
NABirds	SENet-154	88.3	89.2	PC-DenseNet-161	[Dubey <i>et al.</i> ArXiv’17]	82.8
Birdsnap	SENet-154	83.4	84.3	EfficientNet-B7	[Tan & Le, ArXiv’19]	84.3

Table 1: Transfer learning tasks. Evaluation is top-1, single crop and without changing the base Imagenet architecture.