

1 We thank the reviewers for their feedback, and their positive appraisal of our work. We respond to each point below.

2 **R1:** “*I find your experiments extremely interesting with good results and insight, but unfortunately I don’t think the*
3 *current state of the [figure] presentation is publication-quality.*” These fixes are easily made: we will use text rather
4 than WordNet IDs in Fig. 1, increase the font sizes, and add bolding to Tables 2-4 to bring out the key results. We have
5 updated Fig. 2, reducing the widths of the error bars and grouping results by dataset rather than by model architecture,
6 and Fig. 3, adding a color gradient to our markers to indicate the additive accumulation of augmentations. We retained
7 marker shapes for accessibility of colorblind readers and to ensure unambiguous mapping to the legend.

8 **R1:** “*In Figure 5, I would like to see intermediate (not after the final conv) layers included.*” **R3:** “*...why not report*
9 *Figure 5 for every single representations one can get...*” We initially focused on the output of the convolutional layers
10 since this is the convolutional representation most commonly used for transfer to other computer vision tasks and
11 most correlated with IT activity in the primate visual system. We have extended our analysis to consider each layer of
12 AlexNet. The results recapitulate the ones reported: shape is persistently more decodable through the convolutional
13 layers than is texture, and its decodability varies less than that of texture, which rises through them. In the FC layers,
14 shape decodability decreases markedly whereas texture increases.

15 **R1:** “*Why did you opt for Exemplar as one of your methods in S6?*” Our choice of models was motivated by a desire
16 for diversity of approaches that perform reasonably well on ImageNet, as well as desire to connect to the existing recent
17 literature. The Rotation and Exemplar models were recently compared in Kolesnikov et al. 2019. Although there are
18 differences in implementation, Exemplar, Instance Discrimination, and SimCLR are all contrastive learning methods
19 and might be expected to behave similarly modulo augmentation. We are unaware of recent work achieving success
20 with ordinary GANs or VAEs for unsupervised representation learning on ImageNet.

21 **R2:** “*On the broader impact section: I miss any discussion of the negative impact this work can have or rather that*
22 *some of the main findings imply.*” We have added the following text: “People...often have a mental model of computer
23 vision models as similar to human vision. Our findings contribute to a body of work showing that this view is actually
24 far from correct, especially for ImageNet, one of the datasets most commonly used to train and evaluate models.
25 Divergences between human and machine vision of the kind we study could cause users to make significant errors in
26 anticipating and reasoning about the behavior of computer vision systems...At the same time, we recognize the possible
27 negative consequences of blindly constraining models’ judgments to agree with people’s: human visual judgments
28 display many forms of bias that should certainly be kept out of computer models.” Per **R3**’s pointer, we have removed
29 the numbering before “Broader Impact.”

30 **R3:** “*what does ‘remove’ mean exactly?*” In Figure 5, we presented results showing that shape information is down-
31 weighted by (becomes less decodable through) the fully connected layers of the AlexNet classifier, whereas texture
32 remains relatively constant. Our statement that “these models’ classification layers remove shape information” refers to
33 the fact that shape decodability is higher at the final convolutional layer (pool3) than in the fully-connected layers, and
34 for each fully-connected layer, shape decodability is higher for its input than its output. We will add clarifying text.

35 **R3:** “*In the abstract, the paper says the differences ‘arise not from differences in their internal workings, but from*
36 *differences in the data that they see’...yet in the experiments, the authors demonstrate that, with more carefully designed*
37 *regularizations...the model can be pushed to focus more on the shape.*” Our goal with this statement was to emphasize
38 the main finding of our study: while factors such as loss function and model architecture do indeed influence the level
39 of texture bias in a model, they do so to a much lesser degree than do features of the data, e.g. data augmentation.
40 However, we agree that the statement is worded too strongly. We have added a qualifier: “may not arise primarily from.”

41 *Comparison with other work.* **R1:** “*I’d be interested in a comparison of your architecture findings to those of Ilyas et al.*”
42 Ilyas’ main architecture analysis is Fig. 3, which investigates the transferability of adversarial examples from ResNet-50
43 to other architectures. While differences in the set of models studied unfortunately prevents a complete comparison, we
44 do observe a qualitative similarity to their results: just as they find ResNet-50 adversarial examples to transfer better to
45 DenseNet than Inception-v3, we find the shape and texture preferences of ResNet-50 to be more similar to the DenseNet
46 variants that we study than to Inception-v3. **R3:** “*I wonder if this main arguments of [Wang et al. (2020)] also follows*
47 *this understanding of shape vs. texture, why or why not.*” This is an interesting question. Several studies have found
48 that CNNs are sensitive to the Fourier statistics of training data and can use high-frequency features imperceptible to
49 people (also Yin et al. 2019). We see this as related to texture bias, and in our data augmentation experiments, find that
50 Gaussian blur, which removes high-frequency information, reduces texture bias. At the same time, we do not see shape
51 and texture as entirely reducible to spatial frequency information (see e.g. Portilla and Simoncelli 2000).

52 *Textual edits.* Thank you to **R1** for noticing several errors, which we have corrected. *Future follow-up questions.* **R2**
53 notes that an interesting future follow-up experiment would be to fine-tune the self-supervised models using the same
54 data augmentations used in the cumulative experiment. We agree – thank you for this suggestion! We also agree with
55 **R4** that it would be interesting to investigate texture bias in object detection or segmentation models.