

1 We thank the reviewers for their insightful comments. We address individual notes below.

2 **Reviewer 1, Q3:** We agree that the utility of a perceptual metric is much increased if it is made available for public
3 scrutiny and further research. We are committed to making the code and model parameters available online upon
4 publication.

5 **Reviewer 1, Q7:** We have done our best to describe the details of our experiments in the paper and the supplemental
6 material. Due to licensing restrictions, we may not be able to make the training set public, but we will make the code
7 and trained model parameters available to complement the description.

8 **Reviewer 2, Q2:** We thank you for your honest comments. Apparently there is some conceptual confusion here;
9 apologies if this is due to our writing. The novelty of our approach is that we do not rely on supervised training data to
10 solve the task (IQA), which we mention in the title, the abstract, and many points throughout the paper. Instead, we
11 build a representation solely based on combining known physiological properties of the human visual system (scale
12 and shift-equivariance; cell responses in cortical area V1, for instance, have been shown to have these properties) with
13 two theories of the visual system’s high-level computational goals: slowness and efficient coding. As we describe in
14 Sec. 2.1 and 2.3, the efficient coding constraint is explicitly implemented by the compressivity of the IXYZ objective;
15 the slowness principle is implemented by extracting the mutual information of successive frames from the training
16 video dataset, using the objective. The fact that our learned representation is highly predictive of human quality
17 judgements (69.4% accuracy compared to 73.9% agreement from human to human on the 2AFC task [35]) not only
18 represents a contribution in the engineering sense (as a way to reduce costly gathering of human responses compared to
19 supervised learning), but can also be interpreted as evidence that confirms the validity of the slowness and efficient
20 coding principles. We have not focused on the latter point, but we will certainly improve our paper in that regard.

21 **Reviewer 2, Q8:** We are of course aware of [35], as we compare to their empirical results directly in Table 1. Although
22 qualitatively comparing IQA metrics is difficult, we make some interesting observations about their approach vs. ours in
23 Section 3.3. In particular, their method appears to inherit certain representational properties of a classifier, like relative
24 equivariance to substantial amounts of fog, as it is built on pre-trained classifiers, while ours does not. Of course,
25 humans also are sensitive to the presence of fog and many other differences that classifiers are often explicitly trained to
26 ignore.

27 **Reviewer 2, Q11:** We did discuss the broader impact of our work in a dedicated section, as required by the submission
28 guidelines.

29 **Reviewer 3, Q3 & Reviewer 4, Q5:** We agree that more analysis of the learned representation would be helpful, and
30 have already begun preparations to carry out further experiments. We are also interested in evaluating whether our
31 representation could be useful for other tasks than IQA. However, we also think that 69.4% accuracy on the 2AFC task,
32 notably without any additional supervised regression stage, is an extremely strong result compared to 73.9% agreement
33 from human to human [35] (which we will stress in the revised paper). In our view, this result justifies publication.

34 **Reviewer 3, Q3/Q5:** Please note that Section 3.4 is entirely devoted to assessing which of the inductive biases
35 contribute to achieving good prediction results. Notably, a multi-scale architecture, or at least reading off activations at
36 multiple layers, appears to be crucial; compressiveness of the objective appears to improve results substantially. Other
37 architectural choices don’t appear to matter as much. Some desirable ablation experiments are infeasible, e.g., dropping
38 the convolutionality constraint would lead to a many times larger model, making it difficult to train in practice.

39 **Reviewer 3, Q8:** We agree that Figure 3 in particular can have a more descriptive caption, and we will improve this.

40 **Reviewer 4, Q3:** First, let us clarify that the slowness principle, as implemented in our paper, does not necessarily
41 imply temporal persistence of the metric, it implies temporal persistence of the representation in a probabilistic sense
42 (i.e., on average, neighboring frames will have a small perceptual distance, but individual frames may have a large
43 distance under the metric). An intuitive depiction of the concept can be found in [http://www.scholarpedia.org/
44 article/Slow_feature_analysis](http://www.scholarpedia.org/article/Slow_feature_analysis), Fig. 2. It is correct that to train our model, we would ideally use uncompressed
45 video, since otherwise our representation may become insensitive to compression artifacts. In practice, uncompressed
46 video with sufficient variability is difficult to come by. We therefore applied pre-processing to our training set, including
47 downsampling by a large randomized factor, in order to minimize compression artifacts, as described in the appendix.
48 Spot checking the sensitivity of the metric to various corruptions in Section 3.3 and the appendix, we find that it does
49 not have any obvious defects that may stem from compressed training data.

50 **Reviewer 4, Q6:** Thank you for this reference, which we weren’t aware of. We will discuss this paper in the revision.
51 While this work uses deep learning as well, it is an example of a supervisedly trained no-reference metric, which is
52 targeted specifically at predicting quality of images “in the wild” without comparing to an original. Our work presents
53 a full-reference metric and aims more at the question of which inductive biases can help unsupervised learning, in
54 particular to make training perceptual metrics more data efficient.