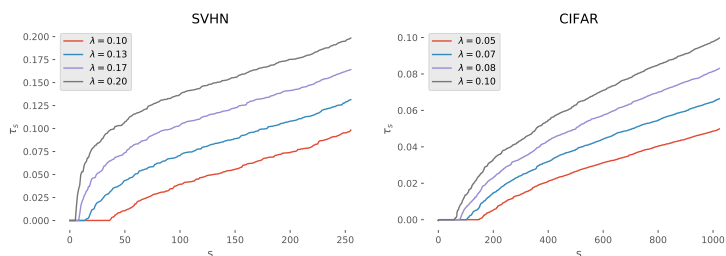1   We thank the reviewers for their encouraging and instructive comments, and the AC for guiding the review process.

2   **Encoder Gap and Further Numerical Evidence (R1, R2 and R3)** The encoder gap can be viewed simply as a measure
3   of maximal energy along any dictionary atom that is not in the support of an input vector. For example, if we assume that
4   the dictionary forms an orthogonal basis for principal subspace of the data, and the encoder is linear (e.g., as in PCA),
5   then the expected encoder gap would be $\lambda - (s+1)^{th}$ largest eigenvalue of the covariance matrix $\mathbf{C_{xx}} = \mathbf{E}[\mathbf{xx}^\top]$
6   (assuming zero-mean data). More generally, it is the $(s+1)^{th}$ entry of the vector $\lambda\mathbf{1} - |\langle \mathbf{D}, \mathbf{x} - \mathbf{D}\varphi_{\mathbf{D}}(\mathbf{x})\rangle|$ (ordered in
7   increasing manner) as we state on line 123. The formal max-min definition follows the previous work of Mehta and
8   Gray (2013), and may look a bit too complicated. We will add a remark in line with our comment above.

9   Note that the assumption on encoder gap is very mild. Intuitively, if a dictionary $\mathbf{D}$ provides quickly decaying
10  approximation error as a function of the cardinality, then a positive encoder gap exists for some $s$. Importantly, the
11  cardinality $s$ provides a knob for our results as one can always consider larger $s$ to guarantee that $\tau_s(\mathbf{x}) > 0$, at the
12  expense of the scaling of our generalization bound (through $\eta_s$). Moreover, one can still induce a larger encoder gap by
13  increasing the regularization parameter $\lambda$, as demonstrated in Fig. 1 in the paper (and the figures in this document),
14  which will come at the expense of accuracy (as demonstrated in Fig. 2c and 2d in the manuscript). In this way, our
15  results guarantee that if one can achieve good accuracy with large encoder gap, then one can generalize robustly. This is
16  reminiscent of generalization bounds for any margin based predictor (e.g., SVM): If the empirical margin loss is small
17  and the margin achieved is large, then the hypothesis generalizes well.

18  Lastly, while our contribution is mostly the-
19  oretical, we also provide further numerical
20  evidence for the encoder gap in real scenar-
21  ios. We trained two models, on SVHN (with
22  256 atoms) and CIFAR (with 1024 atoms),
23  with the training procedure described in our
24  paper, and depict the value of $\tau_s$ (obtained
25  over a collection of samples), for different
26  values of $\lambda$. As you can see, one can easily
27  obtain $\tau_s > 0$ for quite small values of $s$.



28  **R2:** *It is not clear that sparsity-promoting encoders are the right models to be studying.* Deploying sparse priors in
29  the learned representations is a first-take at the analysis of non-linear and data-dependent mappings for supervised
30  models. Ours is the first work to address this. Moreover, *parsimony* (e.g., sparse feature learning) plays an important
31  role throughout data science and machine learning. Sparsity of learned representation can be ensured by sparsity on the
32  weight vectors in the second layer of a neural network (the dense first layer can be viewed as learning a dictionary), and
33  it is indeed common in practice to have an $\ell_1$ penalty on weight matrices. Convolutional structures further promote
34  sparsity. The challenge in analyzing a two hidden layer neural network (as described above) is that it is unclear what a
35  good distributional property would be (akin to encoder gap) that will allow a trade-off between robustness and accuracy.

36  **R2:** *For Theorem 4.1, it would be good to explain how your result compares to prior theoretical bounds on adver-*
37  *sarially robust generalization.* There is no direct analytical comparison since the nature of work here is quite different:
38  none of the prior works study the role of margin and stability of the learned representation in enabling a trade-off
39  between robustness and accuracy. We will add a qualitative comparison and discussion on these in the revised version.

40  **R3:** *It seems that the method is limited to the linear case as well as the incoherent assumption. It is possible that*
41  *the data does not satisfy the incoherent assumption or the data could not be represented by linear combinations of*
42  *columns in D?* Note that the end-to-end map $f(\mathbf{x})$ is *non-linear* in $\mathbf{x}$, as it is the composition of a linear function and a
43  non-linear representation map. Next, the only assumption we need on the incoherence of the dictionary is that $\eta_s < 1$,
44  which is mild. The fact that image data can be sparsified by incoherent dictionaries is well-known (e.g. the cornerstone
45  of JPEG-2000). In practice, this is ensured during training via regularization as described in Equation (7). If you look
46  at the statement of Theorem 4.1, there is no assumption on $\eta_s$, instead our bounds are in terms of $\eta_s$, so the sample
47  complexity is expressed directly in terms of this quantity.

48  **R3:** *It is weird that in Figure 2(b), the accuracy is not monotonically decreasing with regard to the adversarial*
49  *budget.* Computing adversarial perturbations requires solving a non-convex optimization problem. Since this is
50  infeasible in general (for non-linear models), one resorts to approximations (such as those based on projected gradient
51  descent [Madry 2018]). These approximations are not guaranteed to recover the perturbations that maximize the
52  error, and as a result, the empirical reported accuracy is not necessarily monotonically decreasing. Note though that
53  fluctuations are very small and within the margin of approximation.

54  **Experimental details:** We are committed to the reproducibility of our results. All code to reproduce experiments will
55  be shared and openly available.