

1 We thank all reviewers for their constructive comments.

## 2 To Reviewer 1

3 **Q1.** *Other inductive biases likely have an effect:* We agree it is possible and will revise Sec. 4 to better reflect this.

4 **Q2.** *References on non-generative approaches to outlier detection:* Thanks, we will include them in a revision.

## 5 To Reviewer 2

6 **Q1.** *Image data are not exactly temporal sequences:* Even though the data is not truly temporal, the conditional  
7 distributions  $\{p(x_i|x_{<i})\}$  are still well-defined, since the joint distribution  $p(\{x_i\}_{i=1}^d)$  is defined. Thus for a random  
8 variable (rv) in  $\mathbb{R}^d$  following a certain joint distribution (e.g.  $p_{inlier}$ ), we can always view it as a sequence of  $d$  rvs,  
9 where the  $i$ -th rv is sampled from  $p(x_i|x_{<i})$ . Subsequently, we can reason about properties of the random sequence,  
10 such as IID, MD or WN. We will revise the text to clarify this.

11 **Q2.** *How the typicality test is related to IID, WN and MD:* The typicality test is only effective on factorized inlier  
12 distributions. So to apply this test, we need to transform the input rv so that the corresponding distribution becomes  
13 factorized.<sup>1</sup> Usually the transformation is designed s.t. the resulting rv is also componentwise IID. Such a rv in  $\mathbb{R}^d$  can  
14 be viewed as a *sequence* of  $d$  IID rvs. Thus the typicality test amounts to testing the IID condition on the transformed  
15 sequence  $T(x)$ , using the statistic  $\frac{1}{d} \sum_{i=1}^d T_i^2(x)$ . It is not related to MD or WN, which are utilized by our test.

16 **Q3.** *For general AR models, why do we expect the sequence  $R(x)$  to be WN for in-distribution  $x$  and not WN for OOD*  
17  *$x$ :* (i) We showed that  $R(x)$  is always WN when  $x \sim p_{inlier}$  in L101-102. Note that we follow the convention in  
18 the time series literature, and define WN sequences as uncorrelated sequences with zero mean and unit variance; this  
19 definition does not require the sequence to be IID, and is called “weak WN” in some fields. (ii) For OOD  $x$ ,  $T(x)$  not  
20 being WN is the alternative hypothesis we are interested in. This is because designing a universally effective test is very  
21 difficult, if not impossible, given the high dimensionality compared to the limited number of samples. Thus we have to  
22 restrict our attention to certain alternative hypotheses (i.e. certain kind of outliers). In this work, we are interested in a  
23 variety of natural image outliers which have previously led to confusions. We believe that our alternative hypothesis  
24 suits this purpose, as discussed in L121-123 and verified in Appendix A.

25 **Q4.** *Test datasets from other domains:* In this work we are mainly interested in natural image outliers, which have  
26 previously caused confusion over the calibration of DGMs in literature (also see L21-25, L31-36 in our submission).

## 27 To Reviewer 3

28 **Q1.** *Differences between the proposed test and the typicality test, and implications in practice (e.g. what happens if we*  
29 *apply the typicality test on non-IID sequences):* (i) They have different assumptions (WN vs IID), and the test statistics  
30 are different. (ii) The difference in assumptions means that it can be more difficult to apply the typicality test in a  
31 principled way, see L77-80. Even if we consider a heuristic application of the typicality test, i.e. to use the test statistics  
32  $\frac{1}{d} \sum_i T_i^2(x)$  where  $T(x)$  isn’t necessarily IID for inlier  $x$ , the difference in test statistics means that our test can still be  
33 more effective, since it could identify anomalous autocorrelation structures in the (transformed) *outlier* distribution. For  
34 example, suppose  $T_i(x)$  is WN for inlier  $x$ , while for outlier  $x$ ,  $(T_1(x), T_2(x))$  are uniformly sampled from a centered  
35 circle with radius  $\sqrt{2}$ , and for  $i > 2$ ,  $T_i(x) = T_{i-2}(x)$ . Then for such outliers,  $\frac{1}{d} \sum_i T_i^2(x) = 1$ , and the typicality  
36 test will not be able to detect them; in contrast, our test will detect such outliers, since they have *autocorrelations*; see  
37 Remark 2.1. The autocorrelation issue is relevant for natural image outliers, as discussed in L121-123 and Appendix A.

38 **Q2.** *Results of the typicality test when applied on residuals, or the transformed latents of DGMs:* (i) Nalisnick et  
39 al (2019) tested it using a flow model (Fig 4(a) therein) and showed that it was not effective. (ii) Some methods in  
40 Sec 3.1 are equivalent to the typicality test on certain whitened residuals: the LH-2S test using linear models, and the  
41 LN-LH2S method in Appendix (Table 5) which works on VAE residuals. Both methods are clearly outperformed by the  
42 corresponding WN tests. (iii) We also experimented with the typicality test applied to AR-DGM residuals, using the  
43 setup of Sec 3.1. The results are similar, with the typicality test outperformed by ours in 5 out of 6 cases.

## 44 To Reviewer 4

45 **Q1.** *Sample size and CI in Table 1:* For all methods we use the entire test set, except for AR-DGM for which we sample  
46  $5 \times 10^4$  images from the larger datasets. This leads to a minimum sample size of  $10^4$  (CIFAR-10 test set). Using the  
47 formula R4 provided, we can show that the maximum possible 95% CI is  $\pm 0.011$ . We will include them in revision.

48 **Q2.** *How much does the result change using different AR models:* While we can’t train new models due to time  
49 constraints, the results in Appendix B are obtained using a smaller-capacity PixelCNN++, and our test still works well.  
50 Also note that our test works with a simple linear AR model. Based on these results, it seems reasonable to expect that  
51 the results will not change qualitatively as we switch to different models.

---

<sup>1</sup>This is the approach described in Sec. 2.1. See also L179-181 and footnote 1 in the submission, which discussed an alternative weak typicality test; that test is known to be ineffective, likely due to the lack of any concentration guarantee.