

1 We thank the reviewers for acknowledging our contributions and for providing valuable feedback.

2 **To R#1: Transformer latency:** For MNIST experiments, TensorFlow reports that the LeNet-5 BNN requires
3 10,906,677 FLOPS and the WGIN-GP generator requires 9,292,938 FLOPS. The additional cost of the rejection
4 loop is then ~ 20.2 million FLOPS. The NVIDIA Titan X (Pascal) is rated at 11.0 TFLOPS, so the latency of
5 rejection is ~ 0.02 milliseconds on our devices.

6 **To R#3: Additional WGIN-GP architecture and concatenation details:** The architecture of the WGIN-GP
7 generator closely follows that described in the WGAN-GP paper, adding conditioning via the concatenation method
8 referenced in Section 2.1. The conditioning image is flattened and concatenated with a noise vector and passed to the
9 standard WGAN-GP generator as input. The critic architecture also follows the WGAN-GP paper. We concatenate the
10 corresponding one-hot class label to the generated image and, similarly to [*1], to the output of each hidden convolution
11 layer. We clarified the concatenation process in our paper and have added the missing hidden-layer citation. We also
12 added detailed architecture diagrams for each network to our appendix. Thank you for pointing this out.

13 **To R#3: Bayesian neural network and rejection function interaction:** We use Monte Carlo sampling to determine
14 the BNN’s predicted class and uncertainty metric. We first sample the model ten times for the given input x_i , effectively
15 ensembling ten different networks. We calculate the mean of the given class probabilities and treat the argmax as the
16 class prediction y'_i . We treat the median of the probabilities for this predicted class as the certainty metric c_i . These two
17 metrics are passed to the rejection function. We did not see a significant difference in WGIN-GP performance when
18 treating the mean as the certainty metric. Alternative approaches may consider the variance in the predicted class across
19 models. Since uncertainty is represented by the difference in prediction across models, a single network prediction does
20 not provide uncertainty information and thus should not be used alone.

21 **To R#3: Baseline models:** We reevaluate the WGIN-GP using a stronger baseline (expanded LeNet-5 with batch-
22 norm, dropout, and more convolution layers/filters; the exact architecture will be included in the paper) while the
23 rejection function and WGIN-GP are unchanged. The GWIN continues to have a positive impact. Results are shown
24 in Table 1. CIFAR10 experiments are in progress.

25 **To R#4: Denoising methodologies:** A key distinction between our work and the papers mentioned by R#4 is that
26 those papers and their related works focus on using GANs to increase the robustness of a classifier *during training* by
27 generating out-of-distribution training data (similar to hallucination methods in the few-shot learning domain) while
28 our method assumes a *fixed, pretrained* classifier and uses generative methods to translate novel, out-of-distribution
29 examples to the confident distribution *during inference*. Since the GWIN framework learns representations that the
30 classifier labels correctly with high confidence, these generative denoising methods can easily be paired with our
31 framework: a classifier is trained using the aforementioned techniques and the GWIN is then used to transform any
32 novel examples that the new classifier is not entirely robust to. Similarly to DefenseGAN, the flexibility and additive
33 nature of our framework means that we can easily build atop these existing denoising methodologies. Since noise
34 only represents a subset of out-of-distribution observations, we cannot rely entirely on denoising techniques to address
35 classifier robustness. GWINs take a step towards a generalizable, principled framework for “rethinking” uncertain
36 examples and leveraging classifier uncertainty. For completeness, we will add a more detailed Related Works section to
37 the appendix describing how these methodologies address out-of-distribution robustness during the training process.

38 **To R#4: False positive rate:** Fig.5 shows the change in rejected sample certainty of the ground truth label for varying
39 rejection thresholds. The WGIN-GP increases confidence in the correct class more often than not. Depending on the
40 classifier’s accuracy on the rejected subset, the threshold can be tuned to maximize the expected number of correct
41 classifications. Typically, a classifier is required to predict so higher accuracy implies fewer misclassifications and
42 therefore fewer false positives. If true rejection is allowed, a model can still use the GWIN to perform a transformation
43 and reject the original observation if the transformed observation is below some arbitrarily high threshold.

Table 1: Test set accuracy for Digits (top) and Fashion (bottom) on rejected observations using GWIN transformation for the given certainty threshold τ averaged over 10 runs. For $\tau = 0$, acc. is 99.2% (Digits) and 91.0% (Fashion).

τ	% Reject	BNN Acc.	BNN+GWIN Acc.	Rejected Acc. Δ	Overall Acc. Δ	% Error Δ
0.70	0.25	46.44 \pm 12.30	59.96 \pm 7.28	13.53 \pm 15.54	0.03 \pm 0.04	-3.42 \pm 4.15
0.80	0.41	51.79 \pm 7.93	66.13 \pm 8.39	14.34 \pm 11.08	0.06 \pm 0.05	-6.39 \pm 4.90
0.95	0.81	53.44 \pm 3.91	68.60 \pm 4.50	15.16 \pm 4.85	0.12 \pm 0.04	-12.64 \pm 3.66
0.99	1.26	60.11 \pm 4.53	67.63 \pm 3.62	7.52 \pm 3.07	0.10 \pm 0.04	-9.77 \pm 3.69
0.70	2.54	45.24 \pm 2.20	57.10 \pm 2.75	11.85 \pm 3.26	0.30 \pm 0.09	-3.17 \pm 0.93
0.80	4.04	46.49 \pm 2.16	57.49 \pm 2.67	11.00 \pm 2.60	0.45 \pm 0.11	-4.64 \pm 1.09
0.95	8.18	52.23 \pm 1.44	59.18 \pm 1.31	6.95 \pm 2.15	0.57 \pm 0.19	-5.99 \pm 1.93
0.99	12.33	57.04 \pm 1.22	61.61 \pm 0.86	4.56 \pm 1.23	0.56 \pm 0.15	-5.88 \pm 1.48

44 [*1] S. Reed, Z. Akata, X. Yan, L. Logeswaran, B. Schiele, H. Lee, arXiv preprint arXiv:1605.05396, 2016.