

1 We thank the reviewers for their enthusiastic feedback and insightful suggestions. We appreciate the recognitions of
2 contribution in the reviews, such as “the methodology is meticulously detailed, and the work’s contributions to research
3 and development of better generative models, an issue of major importance in machine learning,” that HYPE as a
4 “benchmark for generative models can be very useful for the community to enable more consistent evaluation of new
5 methods,” and that there was “[s]olid writing and justification for HYPE as a benchmark.” We also appreciate the
6 questions and concerns raised in the reviews, as well as the requests and opportunity for clarification in our author
7 feedback response. We address them as follows and in our corresponding revision.

8 **[R1, R2] Proposed diversity incorporation.** R1 and R2 ask how diversity may be incorporated into a full ranking of
9 the GANs. As mentioned in our submission and acknowledged by R1 and R2, the current scope of HYPE is limited to
10 perceptual realism, though we agree that diversity measurement is worthwhile. Similar to HUSE (Hashimoto 2019),
11 we suggest that diversity can be computed using the automated recall score (Sajjadi 2018). Recall measures diversity
12 independently from precision ($F_{1/8}$), the corresponding measure of fidelity in the paper. We will include this comment
13 in our paper, citing both works.

14 **[R3] Related Work.** We will update our draft with an extended discussion on related work, stating how HYPE differs,
15 beyond FID, KID, and $F_{1/8}$. Broadly, HYPE measures human perceptual judgments of generative outputs using humans
16 directly in an inexpensive, widely accessible method. This contrasts with Neuroscore (Wang 2019) because it can be
17 widely accessible by researchers rather than depend on humans with EEGs properly worn and measured. Additionally,
18 the automated metrics 1-NN and Kernel MMD, suggested for evaluation by Xu et al., show success only in a limited
19 architectural case (the convolutional space of an ImageNet-pretrained ResNet) and nevertheless indirectly assess human
20 judgment.

21 Specifically, we address limitations of automatic metrics. Prior work has asserted that there exists coarse correlation of
22 human judgment to FID and IS, leading to their widespread adoption. However, both metrics depend on an ImageNet-
23 pretrained Inception v3 Network to calculate statistics on the generated output (for IS) and on the real and generated
24 distributions (for FID). The validity of these metrics when applied to other datasets has been repeatedly called into
25 question (Barratt 2018, Rosca 2017, Borji 2018, Ravuri 2018). Perturbations imperceptible to humans alter their values,
26 similar to the behavior of adversarial examples (Kurakin 2016). Finally, because FID is a biased estimator, there is
27 inherent variance to the metric depending on the number of images and which images were chosen—in fact, there exists
28 a correlation between accuracy and budget (cost of computation) in improving FID scores, because spending a longer
29 time and thus higher cost on compute will yield better FID scores (Lucic 2018). KID addresses this as an unbiased
30 estimator (Bińkowski 2018), but otherwise strongly correlates with FID and not human judgment, as we report in the
31 paper.

32 **[R1] Impact of number of samples per evaluator.** Similar to our earlier evaluation on increasing the number of
33 evaluators, we find that the CI decreases when increasing the number of samples that each evaluator assesses. Computed
34 via a bootstrap, find that the CI width decreases monotonically when increasing the number of images, specifically
35 from 10.5 to 8.5 when evaluating 20 versus 100 images. We choose 100 images for HYPE infinity because we find it
36 to be a good combination of reliability and efficiency for the GANs we evaluated. This number can be treated as a
37 hyperparameter, and may be increased or decreased based on the desired reliability and quality of GANs evaluated.

38 **[R2] Meaning of HYPE scores.** As R2 points out, HYPE is designed to rank GANs; deltas in HYPE scores may not
39 reflect an equal delta in visually perceived realism. However, we note that HYPE scores are directly correlated with
40 perceived realism, though that correlation may not be linear. Additionally, an absolute HYPE infinity score above 50
41 meaningfully indicates “hyper-realistic” images, or images that appear more realistic than real images.

42 **[R2] Discussion on differences from automated metrics.** R2 requests an extended discussion on results from
43 automated metric differences. We will include comments in the paper on qualitative differences that we believe exist
44 between HYPE scores and {FID, KID, $F_{1/8}$ } scores.

45 **[R2] Dataset and code release.** We plan to release the dataset and code for researchers to analyze after publication.

46 **[R2] Explicit 50:50 real:fake ratio.** R2 mentions that revealing the ratio to evaluators may bias them. We made this
47 design decision intentionally and carefully. Amazon Mechanical Turk forums would enable evaluators to discuss and
48 learn about this distribution over time, thus altering how different evaluators would approach the task (this discussion
49 occurs often). Thus, we decided to make this ratio explicit such that evaluators would have the same prior entering the
50 task. We also make this decision to evaluate a threshold for which fake images appear more realistic than the real ones.

51 **[R2] Hyper-realism with unrealistic real datasets.** R2 mentions the unrealistic appearance of some real datasets such
52 as CIFAR-10. We acknowledge this, and it enables varying levels of difficulty across datasets. Thus, “hyper-realism”
53 is relative to the real dataset on which a model is trained. Some are easier because of lower resolution and/or lower
54 diversity of images.