1 We thank the reviewers for their insightful comments. We first respond to the common concerns across the four reviews
2 before addressing specific reviewer feedback.

3 **Strength of results, hyperparameters, additional experiments:** All reviewers observe we achieve only modest gains
4 compared to the log-uniform spacing in terms of final log likelihood. We emphasize that while this is true, the main
5 benefit of our approach is in avoiding the costly grid required by the log/linear-uniform spacing schedule. To respond to
6 R4's request, as measured in wall-clock time, our GP-bandit schedule takes approx. 12 hrs to train a VAE on MNIST
7 compared to the approx 160hrs training time for the grid searched log schedule (8 hrs/run x 20 runs). We will update
8 the main text to make these considerations clear, and add an additional appendix section benchmarking our schedule
9 against baselines in terms of total wall-clock time.

10 To respond to R4, all results were obtained from
11 a single seed. We will update the figures in the
12 final draft to average results across five seeds, and
13 have included preliminary results on MNIST in
14 Figure 1. As per R1's request, we will also include
15 the training loss in the final version, which closely
16 mirrors the test loss curves.



Figure 1: VAE MNIST, 6k epochs, 5 seeds

17 R2 raises a concern about the number of hyper-
18 parameters in our method. The GP-bandit intro-
19 duces three additional hyperparameters which can
20 be learned directly from data by maximizing the
21 GP marginal likelihood, and therefore require no
22 additional hand-tuning by practitioners as discussed in Appendix B. The exploration-exploitation trade off parameter $\kappa_t$
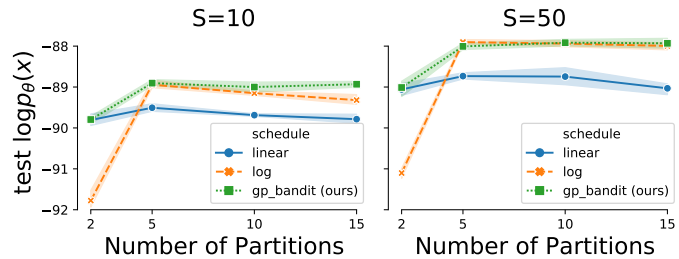23 is set using Theorem 1.

24 **Figure 2 clarity:** We agree with R2 and R4 that a more thorough description of both Fig. 2 and Fig. 3 are needed.
25 In Fig. 2, we investigate the bandit exploration / exploitation behavior at early, middle, and late stages of training by
26 showing where the bandit positions a single $\beta_1$. Color encodes timestep, so clusters of similar colors indicate the bandit
27 is "exploiting" a particular region, particularly in regions where the blue line, i.e. GP mean, is high. This indicates
28 that our reward model expects the TVO objective will improve with this choice of $\beta_1$. We also show the variance of
29 our GP surrogate model across training phases and see both the mean and variance of predicted reward decreases as
30 optimization converges.

31 The key takeaway from this figure is that we see the bandit exploits early on in training, and that the surrogate reward
32 function correctly learns that the reward for choosing the correct location decreases as optimization proceeds and the
33 curve flattens.

34 **Figure 3 clarity:** R2 correctly observes that it is difficult to know the shape of the integrand. We agree, and this is in
35 fact one of our primary motivations for using a bandits schedule, as bandit optimization is uniquely suited to scenarios
36 where one has little knowledge of the form of the reward function. We show possible example shapes for the integrand
37 in the subpanel of Fig. 3, reflecting bandit choices of $\beta$ in the middle (orange) and late (green) stages of training. These
38 are based on the intuition that $\beta$ choices should be concentrated in regions where the integrand is changing quickly,
39 allowing the left Riemann approximation to capture the most area with a fixed budget of $d$ partitions. We also assume
40 that a perfectly flat curve will result in uniform $\beta$ choices. We will update the caption of Fig. 3 to make this clear.

41 **(R4) Comparison with Bogunovic et. al [5]:** We agree that the positioning of our contribution with respect to [5]
42 requires further clarification, and regret this oversight. Theorem 1 improves the bound on the maximum mutual
43 information gain $\tilde{\gamma}$ from [5] by using Cauchy Schwarz and Jensen's inequality rather than an analysis of the optimality
44 conditions (cf. eq. 61 in [5]). We have updated the main text in 4.3 to make it clear that Theorem 1 follows by using
45 our tighter bound on $\tilde{\gamma}$ in Theorem 4 of [5], and simplified the derivation in the appendix to refer to [5] directly where
46 possible.

47 **(R4) Discussion on spatial covariance:** R4 asks for clarification on covariance functions which can be used alongside
48 our projection operation. Since the projection preserves the input space, any PSD covariance function will maintain this
49 property after projection. We recommend choosing a kernel for which bounds on the maximum information gain are
50 known, such as exponentiated quadratic, Mattern, or linear from Srinivas et al 2010 [35].

51 **Typos:** R1, R3, R4 helpfully point out a number of typos and suggestions to improve the writing which we will
52 incorporate into the final draft.

53 **(R4) Ref [43] in the supplement:** Martin J Wainwright. "Basic concentration bounds." In *High-dimensional Statistics:*
54 *A non-asymptotic viewpoint.* Chapter 2, pages 21–57. 2019