1   We thank the reviewers for their useful and thoughtful feedback. We are glad to see that our work was found "*highly*
2 *relevant to the NeurIPS community*" (**R1**) and to be addressing "*a well-motivated problem*" (**R3**) in "*an important area*"
3 (**R1**); that our empirical validation was "*sound*" (**R1**), "*comprehensive*" (**R2**), "*extensive: 7 models are considered and*
4 *two state-of-the-art baselines are included*" (**R3**), and "*thorough* (**R4**)"; with our claims "*supported by a comparatively*
5 *large set of real-world experiments*" (**R1**), results that "*largely align with the claims*" (**R2**) and "*are significant to support*
6 *the claims of the paper*" (**R4**). That is, "*the experiments show that the proposed method outperforms the baselines and*
7 *works consistently well in all the real-world models considered*" (**R3**). We address the reviewers' comments below and
8 incorporated all feedback in the revised paper. Here with 'VBMC' we refer to our method for noisy inference.

9   **Known noise (R1, R4).** The assumption that log-likelihood noise $\sigma_{\text{obs}}$ is (approximately) known is less limiting than it
10 might seem. The common *synthetic likelihood* approach already needs multiple samples to build the multivariate normal
11 pdf used to estimate the likelihood, hence $\sigma_{\text{obs}}$ can be easily estimated via bootstrap with negligible cost (as shown in
12 [10]). *Inverse binomial sampling* (IBS, [41]), used extensively in our analyses, automatically provides both an estimator
13 of the log-likelihood and of $\sigma_{\text{obs}}$ for free. For other techniques, bootstrap is often a feasible and easy-to-implement
14 choice. **R4:** "*How does the estimation error* [of $\sigma_{\text{obs}}$] *affect the results? Will it change the claims of the paper?*" Great
15 question. We performed a new analysis by rerunning VBMC on several problems (aDDM, Timing, Neuronal) while
16 drawing the estimated $\widehat{\sigma}_{\text{obs}} \sim \text{Lognormal}\left(\ln \sigma_{\text{obs}}, \sigma_\sigma^2\right)$ for increasing values of noise-of-estimating-noise, $\sigma_\sigma \geq 0$. At
17 worst, the performance of VBMC via the MMTV metric degrades only by $\sim 0.03$ points on average (e.g., from $0.13$ to
18 $0.16$ on the Timing problem; see Fig. 3 for reference), with $\sigma_\sigma$ up to $0.4$ (i.e., $\widehat{\sigma}_{\text{obs}}$ roughly between $0.5 - 2.2$ times the
19 true value); showing that VBMC is very robust to imprecise estimates of the noise. Thus, our claims are unaffected.

20   **Dimensions & Applicability (R1, R4).** The limit of VBMC to $D \sim 10$ input dimensions is common to approaches
21 that use GP surrogates (e.g., Bayesian optimization [22]), and an open area of research which we intend to pursue. Still,
22 *plenty* of models in computational biology and neuroscience have up to $\sim 10$ parameters, so there is *wide applicability*
23 of our method. As proof, consider the widespread usage within computational neuroscience and related fields of the
24 recent BADS toolbox [26] for (noisy) Bayesian optimization, which shares similar limitations as VBMC in terms of $D$.
25 Crucially, VBMC is a large step forward with respect to BADS for practitioners in the field in that it affords *full Bayesian*
26 *inference* (as opposed to limited to point estimation), so we expect it to impact a wide audience. Finally, while our
27 method does not directly apply to high-$D$ machine learning models (e.g., Bayesian neural networks), it could be used to
28 infer posteriors over *hyperparameters* of ML models, which we see as a relevant research direction.

29   **Theoretical guarantees (R2, R4).** While we share the reviewers' desire for convergence guarantees, we also note that
30 convergence proofs of adaptive Bayesian quadrature methods are outstanding theoretical contributions in themselves.
31 For example, despite a growing body of work on adaptive methods over the years, *only last year* Kanagawa and Hennig
32 [37] were first able to prove convergence for a class of *local* acquisition functions in Bayesian quadrature. By contrast,
33 our paper adds to the literature of solid empirical contributions with theoretically-motivated choices. We believe that our
34 strong empirical validation, judged very positively by all reviewers, while not at all replacing a mathematical derivation,
35 should provide confidence in our method and inspire future theoretical research on VBMC.

36   **Motivation of $a_{\text{VIQR}}$ acquisition function (R2, R3).** The rationale for going from $a_{\text{IMIQR}}$ (Eq. 7) to $a_{\text{VIQR}}$ (Eq. 8) is:
37 (1) sampling from the posterior $\exp(\overline{f}(\boldsymbol{\theta}))$ is relatively hard, whereas sampling from $q_\phi(\boldsymbol{\theta})$ (the variational posterior) is
38 trivial; (2) $q_\phi(\boldsymbol{\theta})$, by construction, approximates $\exp(\overline{f}(\boldsymbol{\theta}))$ up to a normalizing constant (irrelevant for optimization).
39 Thus, we go from Eq. 7 to Eq. 8 by swapping $\exp(\overline{f}(\boldsymbol{\theta}))$ with $q_\phi(\boldsymbol{\theta})$, with substantial gains (see Table 1). While $a_{\text{VIQR}}$
40 adds a layer of variational approximation ($q_\phi(\boldsymbol{\theta}) \approx \exp(\overline{f}(\boldsymbol{\theta}))$), hence the name, it still directly approximates $a_{\text{IMIQR}}$.

41   **Response to remaining comments.** **R1:** *Modelling* $\sigma_{\text{obs}}$, while possible, is tricky in practice due to the trade-off
42 between $\sigma_{\text{obs}}$ and variability of the latent function, which may be hard to disambiguate. However, as argued above,
43 obtaining (approximate) *estimates* of $\sigma_{\text{obs}}$ is very often a viable solution; we leave modeling of *unknown* $\sigma_{\text{obs}}$ as future
44 work. **R2:** In terms of novelty, we propose *variational whitening* (which improves performance on hard problems); we
45 extend (in a fairly straightforward way) previous work from [10,11] to the VBMC framework; we build a novel, extensive
46 noisy benchmark with many models and real datasets from computational and cognitive neuroscience, showing that
47 VBMC vastly outperforms state-of-the-art, providing a meaningful contribution to the field. **R3:** We discussed potential
48 reasons of failure for WSABI and GP-IMIQR in the Supplement, D.2. We believe we performed a fair comparison, in that
49 for both WSABI and GP-IMIQR we tried to fix potential issues, and report the best-performing variants. Re. multimodal
50 posteriors, VBMC (similarly to MCMC) would have trouble with far, disconnected modes, and it is an active research
51 area; the exploratory tendency of GP-IMIQR might help here, but also lead to instabilities. Re. the noiseless case, see Fig.
52 S6, in which we varied the amount of noise; interestingly, most acquisition functions perform similarly with no noise
53 ($\sigma_{\text{obs}} \approx 0$). In Eq. 7, the prior density is included in $\exp(\overline{f})$ ($f$ models the log-joint). **R4:** "*Lack of ablation study.* [...]"
54 The ablation studies are mentioned in the text, and fully reported in the Supplement (E.3, 'Lesion study'). In particular,
55 removing variational whitening degrades performance on the difficult Neuronal and Rodent problems. However, we
56 found no major differences in performance on the other problems; leading us to claim that the improved performance of
57 VBMC on noisy problems (wrt. the original framework [5]) is mostly due to the new acquisition functions.