

1 We thank all reviewers for taking the time to provide detailed feedback and valuable suggestions for our work. We
2 address the reviewers' detailed comments below:

3 • **Comparison with Dasgupta et al. (2017).** We thank Reviewer # 1 for pointing out this interesting work. Both our
4 model and Dasgupta et al. (2017) use flexible transformations to augment a parametric model, but the two approaches
5 differ in (1) the exact type of transformation used and (2) the inference method for the ensemble parameters ω .
6 Specifically, the transformation in our method $G_{\mathbf{x}} : \Phi(\cdot|\mathbf{x}) \rightarrow F^*(\cdot|\mathbf{x})$ is an operator between distribution functions,
7 while the transformation in Dasgupta et al. (2017) is a diffeomorphism $\Gamma_{\mathbf{x}} : y \rightarrow y'$ acting on the response space
8 $y, y' \in \mathcal{Y}$. As a result, the PDFs for the two models are $f(y|\mathbf{x}, \mu) = \frac{\partial}{\partial y} \Phi(\Gamma_{\mathbf{x}}(y)|\mathbf{x}, \hat{\omega}) = \phi(\Gamma_{\mathbf{x}}(y)|\mathbf{x}, \hat{\omega}) * \gamma_{\mathbf{x}}(y)$
9 (Dasgupta et al), and $f(y|\mathbf{x}, \mu) = \frac{\partial}{\partial y} G_{\mathbf{x}}(\Phi(y|\mathbf{x}, \omega)) = g_{\mathbf{x}}(\Phi(y|\mathbf{x}, \omega)) * \phi(y|\mathbf{x}, \omega)$ (this work), where $g_{\mathbf{x}} = \frac{\partial}{\partial \Phi} G_{\mathbf{x}}$
10 and $\gamma_{\mathbf{x}} = \frac{\partial}{\partial y} \Gamma_{\mathbf{x}}$. As shown, the two models share a similarity in that their PDFs are the product of ϕ (a Gaussian
11 PDF) and the derivative of a transformation function. However, their PDFs' exact expressions are in fact different.
12 Comparing the two models, both models are flexible in learning the data-generating F^* . However the construction
13 in this work is better posed for uncertainty quantification in ensemble learning for two reasons: (1) Interpretability:
14 $G_{\mathbf{x}}$ is a transformation that directly "fixes" the original likelihood Φ , which allow us to diagnose the misspecification
15 in original likelihood by inspecting $G_{\mathbf{x}}[\Phi] - \Phi$. This is not easy to do under the formulation of Dasgupta et al; (2)
16 Estimation quality for the ensemble weights ω . Since Dasgupta et al. (2017) estimates ω under the original parametric
17 model, the resulting estimate $\hat{\omega}$ can be biased if the original likelihood is misspecified. In comparison, our model
18 estimates ω under the flexible BNE likelihood to avoid biased inference. In the final manuscript, we will include above
19 discussion in the Related Work section (at the end of Section A.1).

20 • **Implementation difficulty of CGP and HMC.** Since CGP is simply a GP model with probit-based likelihood
21 penalties, implementing CGP is in fact not difficult. Any GP model can be converted to a CGP model by adding log
22 probit terms to the GP log likelihood (e.g., in TensorFlow Probability, this can be done in one line as `gp_likelihood +`
23 `tfd.Normal.log_cdf(g)`) and then apply MCMC as before. The HMC method proposed in this work in fact does not
24 require hand tuning. It uses an automated adaptive step size scheme that is readily available in TensorFlow Probability
25 (see Appendix Section C.3). Alternatively, one can also use the *No U-Turn Sampler* (NUTS) implemented in Stan.

26 • **Flexibility in the mean function.** As shown in equation (5) of the main text, BNE's mean function consists of the
27 original ensemble, the residual process δ and a bias correction term due to \mathbf{G} . In addition to the base predictors, the
28 flexibility of BNE's mean function is mainly driven by the residual process δ , and domain experts can select a flexible
29 kernel for δ to best approximate the data-generating function of interest. (e.g., a RBF kernel to approximate arbitrary
30 continuous functions over a compact support (Michelli et al., 2006)). In the manuscript, we will include this discussion
31 when first introducing the residual process (line 103-106 of the main text)

32 • **Extension to moderately high dimensional predictors.** The BNE framework can be naturally extended to high
33 dimension by choosing kernel functions for δ and \mathbf{G} that are suitable for high-dimensional problems. Example choices
34 include the additive kernel (Durrande et al., 2011) or (deep) neural network kernel (Bach, 2014; Lee et al., 2017).
35 Alternatively, one could also build variable selection into the model using shrinkage priors such as the Automatic
36 Relevance Determination (ARD), spike-and-slab, or Horseshoe (Bobb et al., 2015; Vo and Pati, 2017). In the final
37 manuscript, we will include the above discussion in the Conclusion and Future Work section.

38 • **In the experiment, BNE is by construction more expressive than BAE.** We thank Reviewer # 3 for highlighting
39 an important part of our experiment design. Indeed, BAE is an ablated version of BNE (i.e. without \mathbf{G}). The goals to
40 include BAE are to see (1) if \mathbf{G} leads to significant improvement in large sample sizes, and (2) if \mathbf{G} severely overfits
41 data and leads to worse performance in small sample sizes. Figures 4-5 suggest that the former is true, but not the latter.

42 • **In Appendix H2, the low RMSEs of BNE and BME are within one another's standard deviations.** We thank
43 Reviewer # 3 for making this observation. Indeed, this result is expected and is consistent with what we observed in the
44 simulation experiment (Figures 4-5): In a small sample (where uncertainty is high), BNE is competitive with BME in
45 prediction performance, while providing several advantages when the goal is uncertainty quantification (i.e., uncertainty
46 decomposition in Figure 5, and model diagnosis in Figure H.3).

47 • **Extend empirical comparison.** Following suggestions by Review # 1 and # 3, in the final manuscript we will extend
48 the empirical comparisons in Figure 4-5 to include two more models: the popular Deep Ensemble (Lakshminarayanan
49 et al., 2017) and the classic `npcdist` from the `np` package. Deep Ensemble and `npcdist` fits a finite Gaussian mixture
50 and a kernel smoother to the data, respectively. In terms of performance, Deep Ensemble is expected to perform
51 similarly to BME which also uses a Gaussian mixture. `npcdist` is expected to perform similarly to BNE in this 1D
52 experiment, however generalizing kernel density estimators to higher dimensions is usually more difficult (Scott, 2015).
53 We feel the two scalable GP approaches cannot be directly compared to BNE since they still rely on a parametric model
54 likelihood from a known distribution family. However, we do note that BNE can be made scalable by estimating \mathbf{G} and
55 δ using the variational inducing point method in Hensman et al. (2015). However, the disadvantage of this approach is
56 that the uncertainty estimate may be inaccurate since the variational family usually does not fully capture the posterior
57 distribution. Research into scalable inference method that provides good uncertainty quantification is an important
58 future direction of this work. We will include the above discussion in the Conclusion and Future Work section.