

1 We thank all the reviewers for their time to provide comments and the positive feedback from reviewer 1 and 4. Also,  
2 we appreciate all reviewers for their correct summarisation of our contributions as (1) Introduce the ice-start problem  
3 motivated by real-world examples and proposed a method to tackle it. (2) Propose a novel inference algorithm by  
4 combining amortized inference of local latent variables and a sampling approach for global weights. (3) Propose novel  
5 element-wise acquisition objectives. In the following, we will address the questions proposed by each reviewer.

6 > Clarify definition on two test problems (Reviewer 3) Thanks for the suggestion. We will focus on the imputation  
7 task and make the task setting more clear in the beginning of section 2 and present the prediction task as an extension.  
8 We believe that it is important to include two test tasks as these two tasks together cover the most typical real-world  
9 problems. This demonstrates the broad applicability of the proposed method. Note that we have indeed included the  
10 detailed procedures of both tasks in appendix B.2 Algorithm 3 (imputation) and 4 (active prediction). In the revised  
11 version, we will also add a problem definition section for these two tasks in main text.

12 > Acronyms (Reviewer 3) We will add the NLL acronym in line 80. For AUIC, we used its full name in line 205. The  
13 details about AUIC are included in algorithm 4 in appendix B.2. In the revised version, we will move it to the problem  
14 definition section in the main text.

15 > BELGAM related work (Reviewer 4) Thanks for pointing to further related work, which will be included in revised  
16 version. However, we would like to highlight our novel inference method for BELGAM. The previous related methods  
17 pose strong restrictions/assumptions of the model. For example, they typically require the forward model to be  
18 conjugated to the prior or use mean-field approximations for the posterior. Our proposed method uses a sampling  
19 approach, which is theoretically guaranteed to give the optimal posterior under mild conditions. In BELGAM section,  
20 we will use simple examples to help explaining the generative model for better contextualization.

21 > Questions about the computation of KL term. (Reviewer 4) The KL term  $\text{KL}(q(\theta|\mathbf{X})||p(\theta))$  is indeed intractable due  
22 to MCMC sampling. However, this term appears in the line 136 and is only used to train the encoder parameter  $\phi$ . Thus,  
23 the gradient is taken w.r.t.  $\phi$ , and the gradient contribution of this term is zero.

24 > Encoder for partial observations (Reviewer 1, 4) In the main text, we only briefly mentioned its structure because it is  
25 not our novel contribution. We will add a detailed introduction of this encoder in the appendix of the revised version.  
26 The term  $x_{i_o}$  represents the observed entries for row  $i$ .  $S_i$  is indeed a matrix. Because for each row  $i$ , we have some  
27 observed features. Each feature has a vector embedding. We concatenate each observed feature value (scalar) with its  
28 embedding (vector) to form a vectorised representation of this specific feature. Then, we group all observed features for  
29 row  $i$  and form its matrix representation  $S_i$ .

30 > Notation (Reviewer 4) We will add a table that summarizes the notations used in this paper in the appendix for clarity.

31 > Difference to basic element-wise AL and related work section (Reviewer 1) We will extend and include more details  
32 on the BALD objective in the "Data-wise active learning" subsection, and the comparison to the traditional element-wise  
33 active learning (element-wise AL) in the "Feature-wise active learning" subsection in related work. Briefly, we agree  
34 that ice-start problem is tightly related to the high-level idea of element-wise AL. However we argue that existing  
35 work of element-wise AL methods are commonly restricted to a particular application, such as classification, where a  
36 fully-observed test inputs are required; or associated with strong assumptions, such as linear missing data model (e.g.  
37 matrix completion) or heuristic acquisition objectives. Thus, they cannot handle the problems like active prediction  
38 (test task also has query budget) or imputation tasks with highly non-linear relationship in real-life applications.

39 > Inference algorithm description (Reviewer 1) We thank reviewer 1 for appreciating our contribution of this novel  
40 inference algorithm. We will extend this part, add relevant citations and distinctions compared to previous work in  
41 the revised version. Briefly, there is little previous work that uses Bayesian treatment over weight parameters in the  
42 context of VAE. The previous inference algorithm either relies on the conjugacy of forward and prior distribution, or  
43 strong assumptions over posterior approximations e.g. mean-field. Our proposed method does not requires strong  
44 assumptions over posterior, and is guaranteed to give accurate posterior samples under mild conditions. The better  
45 posterior estimates not only help with the prediction accuracy but will also indeed aid the acquisition because the  
46 acquisition involves the expectation over the posterior. Thus, an inaccurate posterior can result in a poor approximation  
47 of the acquisition and lead to the query of the uninformative feature.

48 > Questions about KL term inside expectation. (Reviewer 1) The reason for writing  $\text{KL}(q(\mathbf{z}_i|\mathbf{x}_i)||p(\mathbf{z}_i))$  inside the  
49 expectation is because we factorize the posterior in line 126 for computational efficiency. Thus, this factorization  
50 decouples the dependency of local latent variables with global weights. In addition, the local variable  $\mathbf{Z}$  is amortized  
51 through each individual  $\mathbf{x}_i$ . The  $\text{KL}(q(\mathbf{Z}|\mathbf{X})||p(\mathbf{Z}))$  can be simply written as a summation of  $\text{KL}(q(\mathbf{z}_i|\mathbf{x}_i)||p(\mathbf{z}_i))$   
52 for all  $\mathbf{x}_i$ . Thus, putting it inside the expectation does not affect its value and we can write it outside in the revised version.

53 > Novelty and combination of acquisition functions (Reviewer 1) To the best of our knowledge, the conditional mutual  
54 information is novel in the context of element-wise AL and the previous work is more based on heuristics (e.g. feature  
55 variance and model improvement). The combination is based on our intuition as described in lines 166-168. But it  
56 indeed corresponds to an information-theoretic quantity (discussed in the paper line 192). This gives an additional  
57 warranty for its validity.