

1 We thank all reviewers for their valuable feedback. Below please find our response to each individual review.

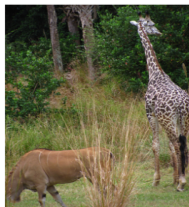
2 **(R1) Significance of improvements:** For VQA, with provided error bars, the improvements are statistically significant.
 3 Especially in the more challenging noisy scenario, the improvements
 4 are over 1 point which is over 10 times the standard deviation. For
 5 NMT, the error bar for BAM-WC is 0.02 (we will add it to the revision),
 6 so the improvement is also statistically significant. Meanwhile, in Table S1,
 7 we compare the run time and number of parameters of BAM and variational
 8 attention [26], where we show that BAM achieves better results while being
 9 more efficient in time and memory.

Table S1: Step time (sec) and number of parameters of variational attention [26] and BAM on NMT.

	S/STEP	PARAMS
VA-ENUM	0.12	64M
VA-SAMPLE	0.15	64M
BAM-WC	0.10	42M

10 **Benefits of modeling attention weights as continuous latent variables:** (a) Modeling attention weights as latent

11 variables enhances the model’s ability to capture complicated dependencies and calibrate
 12 uncertainty, and prevents overfitting due to the added randomness. To evaluate the uncertainty
 13 quantitatively, we provide the PAVPU results for VQA in Table 1 (main paper) and for graph in Table S3.
 14 To evaluate it qualitatively, we visualize the predictions and uncertainties of three VQA examples in
 15 Figure S1. (b) Compared to previous work using discrete latent variables, using continuous
 16 ones is much easier to optimize. Also, BAM is faster and demands less memory and hence more
 17 suitable for low resource scenarios. On hard attention baseline, the choice of REINFORCE
 18 gradient estimator is based on previous work [9, 26]. In [26], Gumbel-softmax, which provides
 19 biased gradients, was found to underperform REINFORCE gradient estimator for NMT, which is
 20 why we have not included it in the experiment.



Question: What animal is next to the giraffe?
 Annotation set: {'wildebeest', 'horse', 'cow', 'antelope', 'gazelle', 'tapir', 'antelope', 'mountain lion', 'antelope', 'horse'}
 Soft answer: deer, p-value: 0.01
BAM-WC answer: cow, p-value: 0.35



Question: What number is on the batter's shirt?
 Annotation set: {'25', '25', '25', '25', '25', '25', '25', '25', '25', '25'}
 Soft answer: 15, p-value: 0.0
BAM-WC answer: 25, p-value: 0.0



Question: Is there mustard on the hot dog?
 Annotation set: {'yes', 'yes', 'yes', 'yes', 'yes', 'yes', 'yes', 'yes', 'yes', 'yes'}
 Soft answer: yes, p-value: 0.48
BAM-WC answer: yes, p-value: 0.0

Figure S1: VQA visualization: we present three image-question pairs along with human annotations. We show the predictions and uncertainty estimates of different methods. We evaluate methods based on their answers and p -values and highlight the better answer in bold (most preferred to least preferred: correct certain > correct uncertain > incorrect uncertain > incorrect certain).

29 **Ablation study on prior distributions:** It is true that the inference network does not get the entire targets during
 30 training, which is the reason that it can be used at the test time to help predict the outputs. We introduce a
 31 contextual prior distribution to impose further regularization on the attention distributions. We agree if setting the
 32 prior and variational posterior the same, the KL in ELBO vanishes and regularization disappears. We have added
 33 an ablation study accordingly, as shown in Table S2, which suggests the importance of appropriate KL regularization.
 34 We will add them into the paper in revision.

Table S2: Accuracy for graphs.

ATTENTION	CORA	CITSEER	PUBMED
GAT	83.00	72.50	77.26
BAM (REMOVE KL)	83.39	72.91	78.50
BAM-WC	83.81	73.52	78.82

36 **(R2)** (1) We clarify that the proposed BAM framework works for any reparameterization distribution defined on the
 37 non-negative real line, and we have chosen Weibull and Lognormal from this family as representative examples. (2)
 38 Compared with existing method, BAM is different in not only variational distribution and gradient estimation, but also
 39 prior distribution. (3) Table S1 shows that BAM is more efficient in both time and memory than variation attention [26].

40 **(R3)** (1) Modeling attention weights is a quite standard approach [9,26,28] as they have intuitive meanings.
 41 Also in BAM, the attention weights are data dependent local variables. This approach is more computationally
 42 efficient compared to the convention in Bayesian neural network where neural network parameters, dependent),
 43 such as θ , are modeled as globally shared random variables (i.e., not data dependent), as the latter approach
 44 needs multiple sampled sets of NN weights to provide uncertainty estimation. In BAM, we only need a single set
 45 of global parameters (NN weights) and rely on the stochasticity on W to provide uncertainty. (2) We did not
 46 choose the gamma distribution as it is not reparameterizable and hence pathwise gradients that are unbiased
 47 and have low variance are not available. (3) We did not include comparison with [26] in VQA as we used a
 48 better baseline model than theirs that had already provided better performance. As their method requires a
 49 case by case design, it is unclear to us how to adapt their method to our model. Further, the code of [26]
 50 was only available for NMT so we only include the comparison for NMT. (4) In Table S3, we have included
 51 the uncertainty estimation result in terms of PAVPU for graph node classification as well and observed
 52 consistent improvements. For other tasks like image captioning, it is unclear to us how to evaluate uncertainty.
 53 (5) In the paper, we eliminated some error bar trying to prevent the table from being too crowded. We
 54 will add them in revision. *Other comments:* In our revision, we will add more detailed explanation for
 55 Figure 1, include the definition of MLE and BLEU metrics, incorporate the missing reference for evaluation
 56 metrics, and update the notation of Equation 1.

Table S3: PAVPU for graphs.

ATTENTION	CORA	CITSEER	PUBMED
GAT	82.30	72.80	77.20
BAM-WC	83.50	73.90	78.10