

1 We would like to thank the reviewers for their feedback and time.

2 **Significance of the work.** R2 and R3 are concerned about the significance of the work. We respectfully disagree with
3 the opinion that the significance is low/mixed. The Bayesian deep learning (BDL) community, which is our target
4 audience, struggles to obtain good performance with VI methods on large problems such as ImageNet (see [3] as an
5 example). This paper is the first to close this gap (as R3 and R5 both mention). Our codebase will enable the community
6 to easily apply VI on large problems, which is a significant contribution. There is a misunderstanding among reviewers
7 regarding our target audience, and we will modify the introduction to make sure that this is fixed.

8 **Significance of experimental results.** R2 and R3 find the experiments not to be convincing and state that our methods
9 do not beat the baselines. We emphasise that this paper is not about beating the state-of-the-art, rather it is about
10 showing that a principled approach works well. The BDL community currently largely relies on MC-dropout, and our
11 goal is to show them that VI can achieve similar or better performance.

12 Our results achieve this objective. Results in Table 1 show that our method performs comparably to SGD, Adam, and
13 MC-dropout. Calibration curves (in Fig. 1) and OOD tests (in Fig. 5) show that uncertainty performance is better than
14 MC-dropout and Adam. It is true that there is no clear winner in Table 1, which is perhaps the reason behind reviewers'
15 concerns. But VOGN does provide a marginally better performance, e.g., on CIFAR-10, on uncertainty metrics, VOGN
16 is consistently either best or tied best (8 out of 12 numbers) or else second best, while both Adam and MC-dropout vary
17 wildly. We will modify the text to make these points clear.

18 **Additional experiments.** R2 and R3 ask for more experiments to
19 show the benefits of Bayesian principles. We will add two experi-
20 ments in the paper (or in Appendices, depending on space constraints):
21 (i) a continual learning experiment, (ii) the Diabetic Retinopathy Di-
22 agnosis benchmark for Bayesian models [5] (this benchmark has
23 only recently been released). For (i), we compare with EWC [4] and
24 Variational Continual Learning (VCL) [1] on 10 tasks of permuted
25 MNIST, using the same architecture and setup as in the VCL paper.
26 Part of the results are shown in Fig. 1 where VOGN achieves $94\pm 1\%$
27 (20 runs), which is better than EWC and SI [1], and marginally better
28 than VCL's performance of $93\pm 1\%$ [2]. An advantage of VOGN over VCL is that it is much faster to converge.

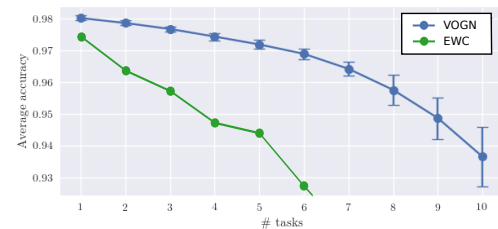


Figure 1: VOGN clearly outperforms EWC (curve from [2]) on 10 tasks of permuted MNIST.

29 **R2: “discuss more about how easy the proposed method can be generalized to different deep models.”** VOGN
30 is a plug-and-play optimiser in PyTorch that is very easy to use (see lines 46-58 in `utils.py`, Supplementary material).
31 Applying VOGN, instead of Adam, requires just 2-3 lines of code change.

32 **R3: “for the out-of-distribution uncertainty ... we don’t know the uncertainty of the true posterior.”** The true
33 posterior is never available in such complex settings, and many other papers have focused on metrics for out-of-
34 distribution uncertainty (see the references in lines 257-258 in the paper).

35 **R3: “... an attempt to gain insight into why VOGN works better than e.g. BBB.”** Thank you for raising this point.
36 As noted by R5, the main reason why VOGN works is the similarity of its updates to Adam, which makes it easier to
37 apply the performance-improvement techniques used in deep-learning. We will add more discussion in the paper so
38 that this point is clearly communicated. We do find that BBB is accurate whenever we can get it to work, but then it is
39 extremely slow to converge; VOGN on the other hand is much quicker (see results on CIFAR-10/LeNet-5). Applying
40 similar techniques for BBB does not work since the updates are very different from Adam (this is what we are saying in
41 lines 91-97: we will modify the text to improve clarity). We have tried many tricks on BBB, including using the local
42 reparameterisation trick and suitably initialising means and variances (as recommended by [2]). We did not specifically
43 try the trick you mentioned.

44 **R3: “I find it surprising that Noisy KFAC is apparently difficult to tune.”** What we meant is that Noisy K-FAC is
45 much slower than VOGN, which makes it difficult to find good hyperparameter settings.

46 **R5: “how your approach compares to stochastic gradient langevin dynamics.”** In VOGN, weights of the neural
47 network are perturbed, while in preconditioned SGLD, gradients are perturbed. If it helps, we can add a simulation
48 comparing the two (although this type of comparison is done previously in [6]).

49 [1] C.V. Nguyen et al. Variational continual learning. ICML, 2018.

50 [2] S. Swaroop et al. Improving and understanding variational continual learning. arXiv:1905.02099, 2019.

51 [3] Y. Ovadia et al. Can You Trust Your Model’s Uncertainty? Evaluating Predictive Uncertainty Under Dataset Shift.
52 arXiv:1906.02530, 2019.

53 [4] J. Kirkpatrick et al. Overcoming catastrophic forgetting in neural networks. PNAS, 2017.

54 [5] A. Filos et al. Benchmarking Bayesian Deep Learning with Diabetic Retinopathy Diagnosis. 2019.

55 [6] Z. Nado et al., Stochastic gradient langevin dynamics that exploit neural network structure. ICLR, 2019.