

1 **R1) Comments on the main proof strategy.** We thank the reviewer for the insightful comments on the proof. We
 2 agree that the sketched argument based on the Bernstein-von Mises theorem is simpler, yet it relies - at least in its
 3 vanilla version- on the assumption that the target function is representable in a fixed parametric model. In order to avoid
 4 this assumption - which is not controllable a priori, we rely on limit results in the over-parameterised regime, in which
 5 the size of the model is larger than the input data. In this scenario, both the input data and the model size are taken to
 6 the limit, and so the limiting model is non-parameteric and guaranteed to have zero loss on the data manifold. It is not
 7 obvious to us if and under which conditions the BvM theorem generalises to this setting, but this is indeed a promising
 8 direction to investigate further. We still believe, though, that our proof provides some interesting geometric insights on
 9 adversarial attacks. We will clarify better in the main text notions like “overparamaterise” or “fully trained”.

10 **R2, R3, R4) Results reproducibility, convolutional layers, and number of different architectures.** We remark that
 11 the hyper-parameters used for training are reported in the Supplementary (Section 3). The source code will be made
 12 available after the review phase. We remark that our experiments comprises both fully connected (up to 5 layers) and
 13 convolutional layers (see Table 2). Noted by R2, Table 4 in the supplementary was accidentally truncated by one value
 14 in several rows, we will update the parameters accordingly (in total we trained 1728 HMC BNNs).

15 **R2) Using deterministic networks and deep ensembles as baseline models.** We agree with the reviewer and in Table
 16 1 we consider the same NN used to perform the experiments in Section 5.2 (hyper-parameters are reported in Table 3 in
 17 the Supplementary) and run a comparison with both deterministic NNs and deep ensembles (Lakshminarayanan, 2017).
 18 We further evaluate the robustness of deep ensembles on a subset of the NNs employed in Section 5.3. We find that
 19 deep ensemble NNs have a robustness similar to that of deterministic NNs suggesting that simply averaging predictions
 20 for different weight initialization and mini-batching is not enough to achieve a robust model. We will add these results
 21 in the main text.

22 **R2) Priors.** In our proof setting, an uninformative prior
 23 is one that gives equal density to all the possible weights
 24 realisations. This can be seen as the limit of a Gaussian
 25 with infinite variance. In practice, the relative importance
 26 of the prior w.r.t. the likelihood diminishes as more data are
 27 used for training, and the posterior distribution gets pulled
 28 apart from the prior. In the experiments reported in the
 29 paper, we have found an $\mathcal{N}(0, 1)$ prior to work well, as the
 30 posterior variance gets to around 0.05 after training. We
 31 performed evaluations with a higher prior variance (up to
 32 10) and noticed a similar behaviour of the loss gradients.

33 **R2) Correlation between accuracy and robustness.**
 34 When a BNNs has a high number of neurons and high
 35 accuracy the conditions for Theorem 1 are approximately met. This guarantees that the network is protected against
 36 gradients attack, thanks to the cancelling effect of Theorem 1. For deterministic NNs Theorem 1 does not hold. The
 37 trade-off between robustness and accuracy in that case has been already observed and studied (Zhang et al.2019).

38 **R3) Proof of Th 1: novelty and completeness.** We would like to stress that Theorem 1 and its proof are novel. In
 39 fact, although we rely on known results for over-parametrised NNs, to the best of our knowledge, the application of
 40 these results in the context of robustness of Bayesian NNs and Lemma 2 are novel. Furthermore, we would like to
 41 clarify that for an input point and a fully trained model, Lemma 2 guarantees that there exists another model such
 42 that the NN has the same loss on the data manifold and opposite orthogonal gradients on that input point. Hence, by
 43 definition, the NN will have same likelihood on models. If they also have same prior (uninformative prior assumption)
 44 then they also have same posterior. This entails the cancellation of orthogonal gradient average.

45 **R3) Are HMC and VI distributed according to the posterior?** Unfortunately, for non-linear networks computation
 46 of the posterior is analytically intractable. Nevertheless, HMC converges to the true posterior in the limit of infinitely
 47 many Monte Carlo samples taken (we used 500 samples in our experiments). On the other side, VI is a more scalable
 48 approximate inference method, but has no convergence guarantees to the true posterior. This also explains why for
 49 MNIST, where both HMC and VI obtain $> 90\%$ accuracy, HMC tends to be more robust than VI.

50 **R4) Why HMC performs better than VI in MNIST but not in FashionMNIST?** HMC converges to the true
 51 posterior, but it is less scalable than VI. As a result, on MNIST where HMC is able to achieve $> 90\%$ accuracy, it tends
 52 to be more robust than VI. On the other hand, on FashionMNIST we were not able to train a BNN with HMC to have
 53 high accuracy, hence, we are far from the regime required by our Theorem 1 to achieve cancelling gradients.

54 **R4) Add comparison with other training methods.** Please, see R2) Using deterministic network and deep ensembles
 55 as baseline models.

Model	Test accuracy	FGSM accuracy	PGD accuracy
Deterministic NN	97.69	21.19	1.45
Ensemble NN	99.4	20.6	0.3
Bayesian NN	96.1	90.0	89.8

Table 1: FGSM and PGD attacks on the network employed in Section 5.2. We compare a deterministic NN, a deep ensemble NN (of size 100), and a BNN (trained with VI). Attacks are performed on 1k test points from the MNIST dataset. We observe that VI trained network achieve better robustness against PGD and FGSM.