1 We are very grateful that all four reviewers recognise the novelty and originality of the paper. Below, we first clarify the
2 main contribution (in particular regarding Lemma 4.3), and then address the detailed points raised by the reviewers.

3 **Main Contribution** is the novel concept weight correlation (WC), which we believe is a key factor affecting the
4 generalisation ability. To consolidate this, we inject WC into the well-received PAC-Bayesian framework to derive a
5 *closed-form expression* of generalization gap bound with *mild assumption* on weight distribution, and then employ WC
6 as an explicit reguraliser to enhance generalisation performance within training. More importantly, the regulariser is
7 effective and computationally efficient in enhancing generalisation performance in practice.

8 Lemma 4.3 shows the positive correlation between WC and the bound (but not the generalisation itself). The effective-
9 ness of WC in either predicting the generalisation (Section 6.1) or training (Section 6.2) is shown with experiments.

10 The central question raised in the reviews is the justification of the assumption that we have a **Gaussian posterior**
11 distribution. We believe that this is a very mild assumption, which is partially justified by distributions tending to
12 converge against Gaussian distributions. More importantly, our assumption significantly relaxes the assumptions used
13 in prior works: an i.i.d. assumption made in [Dziugaite and Roy, 2017, Neyshabur et al. (2017); Jiang et al. (2020)].
14 Different to assuming a general Gaussian distribution, i.i.d. is unrealistic, such that we can lift an unrealistic assumption.

15 This also addresses the **R3: novelty of theoretical claim on WC correlation** concern: we have lifted the—unrealistic—
16 i.i.d. assumption from [Dziugaite and Roy, 2017], landing practical relevance to their findings.

17 This puts us in a sweet spot between the techniques that make unrealistic assumptions about the posterior distribution
18 (usually i.i.d.), and approaches that make no assumptions, but only allow for an a posteriori estimation. (Moreover, such
19 estimations are hard to compute and inaccurate for high dimensional data.) As a result, we have gained the capacity to
20 develop a regulariser, which is both meaningful and easy to compute.

21 **R1, R5: comparison with the state-of-the-art** While mainly focused on [Chatterji et al., ICLR'20], our results also
22 shed light on the fantastic measures paper [Jiang et al, ICLR'20]. In particular, in [Jiang, et al, ICLR'20], it is concluded
23 that "Sharpness-based measures such as sharpness PAC-Bayesian bounds ... perform better overall and seem to be
24 promising candidates for further research". Our results advance this and show that the PAC-Bayesian bounds can be
25 further improved with the weight correlation. Therefore, we believe our paper has advanced the state-of-the-art.

26 **R1, R2, R5: Figure 1** Figure 1 (in Introduction) is an illustrative example to show the positive correlation between WC
27 and generalisability. It is *not* to suggest a general trend for WC (the evolution of WC can be fluctuating) or a conclusive
28 result for the positive correlation. The latter is obtained by the theoretical and empirical results in the following sections.

29 **R2, R5: theoretical support of Lemma 4.3 to regulariser** We agree PAC-Bayes can only provide partial theoretical
30 support to a regulariser. Lemma 4.3 shows that WC is positively correlated with the generalisation bound. This result
31 suggests that it *may* be effective to consider WC as a regulariser. Then, extensive experiments are conducted to validate.

32 **R3: discussions/justification on the particular form of the posterior distributions** We guess you are referring to
33 Def. 4.2. Given the weight matrix $w_\ell \in \mathbb{R}^{N_{\ell-1} \times N_\ell}$ with $i$-th column $w_{\ell i}$ as a random vector, the posterior covariance
34 matrix $\Sigma_{Q_{w_\ell}}$ is defined in a standard way as $\Sigma_{Q_{w_\ell}} = \mathbb{E}[\text{vec}(w_\ell)\text{vec}(w_\ell)^T] \in \mathbb{R}^{N_\ell N_{\ell-1} \times N_\ell N_{\ell-1}}$, where $\text{vec}(\cdot)$ is the
35 vectorisation of a matrix. The $(i,j)$-th block is $[\Sigma_{Q_{w_\ell}}]_{i,j} = \mathbb{E}[w_{\ell i} w_{\ell j}^T] \in \mathbb{R}^{N_{\ell-1} \times N_{\ell-1}}$. For computational simplicity,
36 we use the arithmetic mean instead of the expected value, so that the weight correlation $\rho(w_\ell)$ can be used to represent
37 $[\Sigma_{Q_{w_\ell}}]_{i,j} = \rho(w_\ell)\sigma_\ell^2 I_{N_{\ell-1}}$. Therefore we have $\Sigma_{Q_{w_\ell}} = \Sigma_{\rho(w_\ell)} \otimes \sigma_\ell^2 I_{N_{\ell-1}}$, where $\otimes$ is the Kronecker product.

38 **R3: modest improvement in Section 6.2** While the improvement might look modest in individual cases, it is persistent
39 across the experiments.

40 **R3, R5: more experiments** We conducted additional experiments on Caltech-256 dataset for this rebuttal. It is
41 unrealistic to consider ImageNet, since a single training may take days (or months). The results on the Caltech-256
42 dataset are also promising, and similar to MNIST and CIFAR10. For complexity measure (like Table 2), we have
43 Kendall's $\tau$ at 0.33, as opposed to others at 0.28, 0.28, 0.17, 0.22. For the comparison of models with and without WCD
44 (like Table 3), those models with WCD achieve about 1% improvement on top-5 error over models without WCD.

45 **R5: WCD and posterior** We can confirm that our WCD is obtained from the expression of the posterior of PAC Bayes
46 (Lemma 4.3), and therefore there is no discrepancy between posterior of Lemma 4.3 and WCD-based posterior.

47 **R5: technical details of Lemma 4.3** There may be some misunderstanding to our technical details. For example, the
48 prior is a diagonal matrix $\text{diag}(\sigma_\ell^2)$ and the PCA Bayesian is conducted by summarising the results across layers.

49 **R5: termination condition of training** We use the same termination condition across all our experiments, so it is fair
50 to all methods. The reason for us to consider the best loss across the last few training epoch is to make sure that our
51 results are not affected by random factor – we find that there may still be uncertainty in the last few epochs.