1 We want to express our gratitude to be part of a community that is so generous in offering such insightful and thorough
2 reviews. We were glad that reviewers thought that the paper was "extremely well written" with analysis that is "great
3 and easy to follow" and a "strong step towards seriously understanding why MF assumptions are fine in practice with
4 deep networks" [**R1**], "timely and interesting" [**R2**], with experiments that are "genuinely interesting to think about, as
5 well as genuinely novel" [**R3**], and "significant and novel" [**R4**].

6 [**R1**,**R3**,**R4**] **Motivation:** **R1** felt MFVI already shows "strong performance", excepting optimization problems, and
7 **R4** agreed its use was "common practice". In contrast, **R3** was unsure anyone "would want to use MFVI". 1) The
8 disagreement between reviewers shows how timely and relevant our work is. 2) This highlights how difficult it is to pin
9 down 'what everyone believes'. 3) To substantiate our view of what many published authors believe, our work cites 12
10 papers in §2 that say something like e.g. "a fully factorized posterior... is very restricting" [Louizos and Welling 2016.].
11 4) We agree about optimization issues, and discuss this on L109-110 as well as in S7. 5) We like the reframing **R1**
12 suggested, that our paper justifies common practice in addition to challenging a mistaken belief, and will add this. **But**
13 **whatever the reviewers personally believe, many publish based on the beliefs we critique, motivating our work.**

14 [**R2**,**R3**,**R4**] **Figure 2—in-between uncertainty**: This *toy* fig. shows that in-between uncertainty is possible in deeper
15 networks; larger-scale exps. should carry much more weight. But below we show that for 3-layer MFVI the uncertainty
16 is worse but of the same order-of-magnitude as an HMC baseline, while in 1-layer MFVI fails to capture uncertainty.

17 [**R2**,**R4**] **Special cases:** [**R2**,**R4**] **ReLUs**: Thm. 2 applies to all piecewise linear (PWL) activations (L160). This
18 includes ReLUs and many variants, which are very widely used. Many more theoretical papers focus on PWL, and
19 other nonlinearities can be approximated as PWL. [**R2**] **Gaussians**: in fact, lemmas 1 and 2 and Thm. 2 apply to any
20 distribution with finite first- and second- moments, which we will clarify. Also, fully-factored Gaussians are very
21 commonly used. [**R4**] **MVG**: We give the *general* form of the covariance (L125, D.1 eq. 15) and show equivalence to
22 the commonly used MVG distribution only as an intuitive lower-bound on the expressiveness of the covariance.

23 [**R2**,**R3**] **"Depth induces correlation is well-known"**: **R3** cites Hinton & van Camp 1993 and **R2** offers no citations.
24 H&vC *hypothesise* MF networks could learn complicated posterior distributions over functions, a two sentence comment
25 in their discussion (which we will cite). We *characterise* the form the correspondence takes and compare to the popular
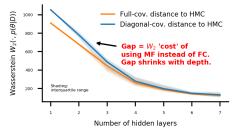26 MVG approximation. Moreover, many papers are motivated by the opposite view, see above.

27 [**R1**] **Completeness of Thm. 2**: This is a subtle point which we will make clearer by making Thm. 2 more mathemati-
28 cally precise. Thm. 2 is an existence proof ("can have non-zero covariance" L185) so proof by construction is valid
29 (L808). But it does not entail the WDH. We then make the admittedly informal argument based on experiment, which
30 you criticize, that the WDH represents a typical case. We can further develop the discussion of this in L189-191.
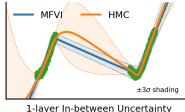
31 [**R1**] **Does full-covariance ever help?**: 1-layer FC beats MF. See below HMC or UCI in Louizos & Welling [2016].
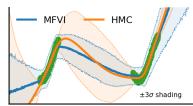
32 [**R2**] **Correlations v. structure**: All exps. compare same model w/ and w/o correlations, which disambiguates this.

33 [**R3**] **MFVI over-confident?**: 1) Thanks for Wasserstein suggestion. See LHS fig. below comparing $W_2$ distance of FC
34 and MF Gaussian to HMC samples. The gap between $W_2$ lines is the distance cost imposed by using MF instead of FC.
35 The gap shrinks with depth, and is negligible after 4 layers, supporting our argument. **This directly addresses your**
36 **point and might justify an increased score**. 2) To clarify: our scope is defending MF assumption specifically, we
37 agree that VI and unimodal apprx. have problems which are not in our scope. 3) Table 2 ECE: As you say, for SWAG,
38 MF ECE is worse, but VOGN v. Noisy K-FAC tells opposite story. This makes sense if some but not all modes are
39 apprx. MF. VI finds the MF modes but SWAG training is indifferent to this (L310-313). Also ECE differences are small.

40 [**R3**] **Non-Bayesian baseline**: Fascinating point: note deterministic baseline has no correlations! Success of determin-
41 istic NNs supports our argument. We will add this. Note also a general defence of Bayesian methods is out of our scope.

42 [**R3**] **Novelty**: Only Thm. 3 draws from Foong et al. 2020, which we are very explicit about, while the rest of our
43 results were developed in parallel. Thm. 3 is only half a page of our paper, which makes it hard to see how this result is
44 "overblown". We will add discussion of Cherief-Abdellatif [2020], but observe (as you highlight) the setting considered
45 has significant differences. We feel that while you acknowledge the novelty of our experimental contribution, you might
46 also ascribe too much prescience to H&vC (see above) and have neglected key contributions (Thm.s 1 and 2).

47 [**R4**] **Exp. critique & Other**: 1) 1-layer HMC is below center. 2) Fig 3 sadly must be small (see C.3) but it does show
48 MF beat FC (maybe unclear in black-and-white, we will make fig. less color-dependent). 3) We discuss KL not going
49 to zero in L268. But also see Wasserstein fig. below LHS, where gap nearly disappears. 4) We agree deciding sufficient
50 width is an unsolved problem and list this and other problems in L214-221. 5) Thanks for the detailed further comments.



Left figure: Wasserstein $W_2(\cdot, p(\theta|D))$ vs Number of hidden layers. Legend: Full-cov. distance to HMC; Diagonal-cov. distance to HMC. Annotation: Gap = $W_2$ 'cost' of using MF instead of FC. Gap shrinks with depth. Shading: interquartile range.

Middle figure: 1-layer In-between Uncertainty. Legend: MFVI, HMC. ±3σ shading

Right figure: 3-layer In-between Uncertainty. Legend: MFVI, HMC. ±3σ shading