

1 We would like to thank all three reviewers for their constructive assessment of our work.

2 **Reviewer 1**

3 *Importance of the experimental results:* Non-vacuous generalization bounds are arguably desirable for delicate machine
4 learning applications. Since such results are actually rare for neural networks, we consider important to show empirically
5 that our approach leads to sound performances. We believe that this can encourage others to develop similar methods
6 that will eventually lead to higher practical impact, following this work.

7 *Error bars:* The reviewer concern is *partially* addressed in Section B.3 of the appendix, through the Monte Carlo
8 sampling size effect study. Indeed, Figure 5 displays the test error — with error bars — for PBGNet and PBGNet $_{\ell}$
9 when varying the sample size. Note that each result is obtained by averaging over 20 repetitions of the learning
10 procedure, each of them executed on different (random) train/test/valid dataset splits, and the stochastic gradient descent
11 is initialized with different random weights. That being said, we undertake to push forward the variance analysis for the
12 final version of the paper, as detailed in Reviewer 3 *Experiments* section.

13 **Reviewer 2**

14 *Choice of the activation function:* A variety of activation functions exists in the literature and we obviously do not focus
15 on the most common one. Nevertheless, the sign activation function is used in the binary networks cited in the paper,
16 notably to reduce the memory footprint of such networks which could be embedded on small devices (e.g., Bengio,
17 2009). In our context, the sign activation is crucial to apply the mathematical trick, and express the predicted outcome
18 as the erf function. This gives a principled way of training the network and derive the PAC-Bayes generalization bound.

19 *Choice of prior:* The bound holds regardless of the choice of the prior μ . In our experiments, we centered the prior on
20 the SGD initialisation weights (as in Dziugaite and Roy, 2017), which corresponds to the real-life scenario where one
21 does not have prior knowledge about the task at hand. The PBGNet $_{pre}$ variant of our algorithm opens the way for using
22 a prior from a precedent learning task, as the transfer learning scenario mentioned in the paper (see Lines 219-222).

23 *Mismatch between theory and experiments:* We worked hard to provide a rigorous and honest empirical study of our
24 theoretical analysis strengths and weaknesses. It is certainly disappointing that there exist good predictors with trivial
25 bounds (as mentioned by Rev. 3), but we still managed to obtain very tight bounds for more than decent predictors.
26 Note also that we compared to tanh networks as this activation is similar to the erf function derived from our analysis.
27 It allows us to use the same optimisation scheme and hyperparameter grid to compare the methods on equivalent basis.

28 **Reviewer 3**

29 *Improvement due to the binary activations:* Relying on binary activation function allows us to express a close-formed
30 solution for the PAC-Bayes bound, without other assumptions than the *iid* one. Another appeal of our bound: it relies
31 more on the network architecture (e.g., d_k , the layer width of each layer k , appears in Eq. 16) than previous results (in
32 the seminal work of Dziugaite and Roy, 2017, a posterior distribution is used over a set of a neural network weights,
33 without taking into account any other architecture specific information).

34 *Bounds for (a unique) binary activated networks:* The reviewer is right to say that our analysis does not provide
35 guarantees for a single deterministic binary activated (BAM) network, but for a continuous aggregation of such BAM
36 networks. We clearly mention this in the introduction (Lines 26-27), but we agree that Line 16 may be ambiguous
37 and we will rephrase it. That being said, several points are worth mentioning: (i) Even if computationally expensive,
38 our predictor closed-form expression is deterministic; (ii) The prediction using Monte Carlo sampling empirically
39 shows a small standard deviation consistently below 10^{-3} , as discussed in the appendix (see Lines 432-448); (iii) In our
40 experiments, we observed that predicting with the single *Maximum-A-Posteriori* BAM network generally gives results
41 remarkably close to the aggregated PBGNet predictor (hinting that the posterior may be quite peaked). Recall that this
42 unique BAM network represents exactly the same prediction function as its mapped tree predictor. This preliminary
43 observation suggests that the bound can provide an appropriate guide to train a BAM network.

44 *Link with similar optimization procedures:* There is definitely a strong connection between our optimization procedure
45 and the REINFORCE method. We greatly thank the reviewer for pointing us the relevant literature; we were genuinely
46 not aware of it. Indeed, we will rewrite parts of Section 3 to highlight these connections, and express our sampling
47 scheme as a particular case of a general technique rather than a new one. As a matter of fact, we think this will enrich the
48 paper, as our PAC-Bayesian analysis support existing strategies to train non-differentiable neural networks. Moreover,
49 the connection with Variational Bayes (for a fixed C) is also very relevant, and will be mentioned in the paper.

50 *Experiments:* We share the reviewer concerns about the lack of standard deviation analysis in Table 1. Given our
51 computing resources, we cannot provide these results in this rebuttal. Nevertheless, we commit to produce for the final
52 version of the paper a thorough stability analysis with 20 different random train/test/valid splits for all six datasets and
53 five considered models of Table 1, for a slightly reduced hyperparameters search grid. Furthermore, we will add a study
54 of the training sample size effect on PBGNet $_{\ell}$ and PBGNet with fixed parameters for the biggest dataset (mnistLH).