1 We thank all the reviewers for their insightful comments, suggestions, and references.

2 **Reviewer 1:**

3 *Novelty of tandem loss:* it is not new, but we were not aware of the prior work, we thank Reviewer 2 for bringing it up.

4 *While most of the computed bounds are non-vacuous, they look to be not that tight. Some discussion of this would be*
5 *valuable. Also a discussion of potential ways to obtain tighter bond values, or whether there is a fundamental limitation.*

6 We provide some discussion in Sections 3.2 and 4.4. The major challenge is the estimation of the tandem loss, which is
7 based on overlaps of OOB samples, which are small [see Section 4.4]. It could be that a better estimation technique
8 could be designed in the future. Another limitation is the oracle bound. For example, in the independent case with the
9 growth of the number of classifiers it converges to $4L(h)^2$, whereas $L(\mathrm{MV}_\rho)$ converges to zero [see Section 3.2].

10 **(*)** *Reproduction of C-bounds:* Reproduction of C-bounds requires definition of a margin. We believe that adding it to
11 the body might divert the attention from the main thread of the paper, but we will be very glad to add a section in the
12 appendix, where we introduce the margin, discuss the relation with the tandem loss, and provide the C-bounds.

13 *Line 471:* Yes, you are right, thank you!

14 **Reviewer 2:**

15 Thank you for providing references to "joint error" and Equations 7 and 8 in Lacasse et al., we will add them.

16 *Reproduction of C-bounds:* see our reply **(*)** to Reviewer 1.

17 **(**)** *Adaptation of the strategy from the "second C-bound form":* Note that in our application the first order loss and the
18 disagreements are estimated on different subsamples. The first order loss is estimated on OOB samples, whereas the
19 disagreements are estimated on overlaps of OOB samples and unlabeled data when available. It is not immediately
20 clear whether joint estimation would give an advantage, we will look at it in future work.

21 **Reviewer 3:**

22 *Posterior optimization*

23 We agree that posterior optimization did not improve the test error. However, we note that in prior work posterior
24 optimization was either impossible (in C-bounds, except for highly limiting cases of aligned posteriors in binary
25 classification) or led to considerable deterioration of the test error (as we demonstrate for the first order bound).
26 Therefore, we see absence of deterioration of the test error as a step forward relative to prior work.

27 *1) compute the multi-class C-bound*

28 The multi-class C-bound based on the $w$-margin with $w = 1/2$ (Corollary 1 of Laviolette et al.) is equivalent to our
29 oracle C-tandem bound in Theorem 6. The values of empirical C-tandem bound are reported in the paper. We note that
30 there may be multiple ways of going from the oracle to an empirical bound, but not all of them are directly applicable in
31 the OOB setting, see our reply **(**)** to Reviewer 2. We also note that the general multi-class C-bound (Theorem 2 of
32 Laviolette et al.) cannot be evaluated of the OOB setting, because the max operator in their definition of the margin in
33 equation (3) cannot be exchanged with expectation and the estimation cannot be done using pairs of hypotheses.

34 *2) Please, re-prove Lemma C.14.*

35 Oh, sorry, we missed that the square was outside the expectation on the left and inside on the right when we canceled
36 the terms. Thanks for catching it! The fix is easy. Instead of canceling, take $-\mathbb{E}\left[X\right]^2 \varepsilon^2$ to the right hand side. Then
37 $\mathbb{E}\left[X\right]^2 \varepsilon^2 - 2\varepsilon \mathbb{E}\left[X\right]\mathbb{E}\left[X^2\right] + \mathbb{E}\left[X^2\right]^2 = (\mathbb{E}\left[X\right]\varepsilon - \mathbb{E}\left[X^2\right])^2 \geq 0.$

38 *3) move the posterior optimization experiment to the supplementary and insert some other results instead (Fig H.15?)*

39 We will have an extra page if accepted, so we can have both in the body.

40 *4) Line 281*

41 The kl-inequality has an upper and a lower inverse, which give an upper and a lower bound, respectively. We have used
42 the lower bound.

43 Thanks a lot for the references to additional work on multiclass classification!

44 **Reviewer 4:**

45 *lines 95, 97, and 98:* You are right, thanks!

46 *line 109:* We will add extra brackets, thanks!