1   We thank the reviewers for their feedback. Below we respond to some of the main concerns.

2   **Clarification of the experimental goals** R1 and R3 were dissatisfied with the small scale of our experiments. We are
3   happy to run any additional experiments that are deemed crucial for better understanding of our method. In fact, we are
4   happy to leave the choice of additional data and ensemble methods up to the reviewers. However, we would like to
5   first clarify what we were hoping to illustrate with the current setup, and further discuss what constitutes meaningful
6   comparisons for a wrapper method such as our J+AB that comes with a model-free guarantee.

7   The primary goal of our experiments is to demonstrate that our method achieves near $1 - \alpha$ coverage numerically
8   (according to the theory, $1 - 2\alpha$ is guaranteed). The secondary goal is to verify that although J+AB runs faster than
9   J+ ENSEMBLE, $\hat{C}_{\alpha,n,B}^{\text{J+AB}} \approx \hat{C}_{\alpha,n,B'}^{\text{J+ ENSEMBLE}}$. The final, lesser goal is to relate known stabilizing properties of bagging by
10   comparing J+AB (or J+ ENSEMBLE) vs J+ NON-ENSEMBLE. These goals are either stated or implied in Lines 251-6,
11   294-303, but we promise to make them more explicit in the camera-ready version.

12   There were two main reasons for running experiments on a small scale with a couple of data sets. The first is the page
13   limit. The second is the cost of running J+ ENSEMBLE; we needed the experimental parameters to be quite small to be
14   certain of obtaining results by the deadline. This is remarked in Lines 284-6. However, given our experimental goals,
15   we did not see the lack of scale as a significant defect. The advantage of model-free framework is that our coverage
16   guarantee is impossible to break irrespective of the data and the choice of ensemble. We can always choose to look
17   at more data sets and more ensemble architectures, but this will only produce more plots that all look very similar.
18   Meanwhile, the secondary goal amounts to a sanity check, and we have said that the final goal is of lesser significance.

19   The issue of width is certainly of interest, but here, we would argue that the only meaningful comparisons are those with
20   other wrapper methods. For example, a split conformal variant is competitive with J+AB in terms of computational
21   cost, but is expected to lose in terms of statistical efficiency. Although this is rather obvious by construction, it may be
22   interesting to investigate whether this would translate to meaningful differences in performance. This is a comparison
23   of efficiency that we are happy to add to our current results. Otherwise, the precision of the intervals would be most
24   heavily affected by the fit of the chosen ensemble with the data. However, as we have ceded this choice to the user,
25   opting to develop a fully flexible method that works irrespective of the quality of this choice, we believe that additional
26   comparisons involving more particular instances of $\mathcal{R}$ or $\varphi$ are not as useful and tangential to the topic.

27   **Breakdown of computational complexity** R3 requested a summary of computational complexity. Here, we provide
28   the total number of occurrences for three different types of operations, which can be used to derive the final cost. We
29   focus on bootstrapping, and match the number of models as in Supplement, Lines 166-9. The table below demonstrates
30   that if the model-fitting cost dominates, the cost of J+AB is roughly that of obtaining a single ensemble prediction. We
31   do not claim any advantage for our method when the cost is dominated by aggregation or evaluation. See Line 170.

32

|  | #calls to $\mathcal{R}$ | #evaluations |
|---|---|---|
| JACKKNIFE | $n + 1$ | $n + 1$ |
| JACKKNIFE+ | $n$ | $2n$ |

|  | #calls to $\mathcal{R}$ | #calls to $\varphi$ | #evaluations |
|---|---|---|---|
| J+ ENSEMBLE | $B'n$ | $n$ | $2B'n$ |
| J+AB (on average) | $B'/(1 - \frac{1}{n})^m$ | $n$ | $2B'/(1 - \frac{1}{n})^m$ |
| ENSEMBLE | $B'$ | $1$ | $B'$ |

33   **Tradeoff between computational and statistical efficiency** R5, as well as R3, expressed concerns about a tradeoff
34   between computational and statistical efficiency for the J+AB vs J+ ENSEMBLE comparison. The short answer is that
35   one method does not always win over the other. See Figure S2. First, since $\tilde{B}$ is a user-specified parameter, it can be
36   picked so that the numbers of models in $\hat{\mu}_{\varphi \setminus i}$ are matched on average. See Supplement, Lines 166-9. Second, the more
37   important difference is the *correlation* among $\{\hat{\mu}_{\varphi \setminus i}\}_{i=1}^{n}$. Conditional on the observed data, $\{\hat{\mu}_{\varphi \setminus i}\}_{i=1}^{n}$ are dependent
38   in the case of J+AB and independent in the case of J+ ENSEMBLE. (Note that unconditionally they are always highly
39   correlated for both.) What this means for the precision is expected to depend on the data and the choice of ensemble. In
40   any case, this difference is expected to be much smaller than, say, that for the J+AB vs split conformal comparison.

41   **R1** 2) The takeaway of Table 1 is in Lines 289-91 (as well as Lines 251-2). It is not our goal to see which among the
42   nine (all *instances* of J+AB) performs best. 3) The results for RF vs RIDGE are completely expected given the known
43   results on bagging. The point we are trying to illustrate is in Lines 295-7. 4) Figure 1 shall be amended.

44   **R3** 3) The number of trials was *doubled* for Supplement C, which reduced the standard errors. 4) J+AB can be applied
45   to any ensemble algorithm in the form of Algorithm 1 as long as it is agnostic to the ordering of the input data. $\mathcal{R}$ or $\varphi$
46   may be arbitrarily complicated, e.g., $\mathcal{R}$ may involve built-in hyperparameter tuning; $\varphi$ may have adaptive weights. 1)
47   No further condition on $\tilde{B}$ is necessary. The validity comes from construction of an exchangeable array of residuals. 2)
48   For fully distribution-free guarantee, a random $B$ is necessary, as the array is not exchangeable with $B$ fixed. Figure S3
49   assumes *bagging*, so fixes a particular resampling and a particular $\varphi$. Also, the title and the abstract will be amended.