1 We thank all the Reviewers for their feedback and their service to the community. We are glad that you have
2 all understood the relevance of our work and, in general, appreciated it. In the following, we try to comment
3 on all the raised issues, and how we will address them, which will certainly improve our work.

4 **Localization** We agree that the relationship between empirical centralization and localization is important
5 and we gave it much thought. We explained this relationship in an earlier submission, but had to omit it in
6 the NeurIPS submission due to space limitations. We summarize it in the following paragraphs.

7 Empirical centralization is *complementary and orthogonal with localization*, as it fixes a different issue:
8 localization is akin to a second-moment normalization (i.e., dividing out variance), whereas centralization is a
9 (complementary) first-moment normalization technique. Because both localized and centralized Rademacher
10 averages are themselves Rademacher averages, they are mutually compatible. We recommend centralizing
11 before localizing, as this approach addresses several issues with standard localization methods.

12 In particular, taking localization to mean "analysis of the variance-constrained star-convex hull", function
13 families containing $\{x \mapsto 0, x \mapsto 1\}$ (for example, any case of realizable symmetric classification with 0–1 loss
14 must include these two functions) suffer from $\Omega(1/\sqrt{m})$ convergence rates, since constant functions have
15 variance 0, whereas with centralization, these two functions are the same, removing this lower bound.

16 Instead taking localization to mean "constrain by raw-variance," the above issue is solved, but now the
17 localized Rademacher average bounds are $\Omega(V/\sqrt{m})$, where $V$ is still only the smallest *raw*-variance.

18 By centralizing before localizing, the difference vanishes, as with mean 0, raw and centralized variances
19 coincide. Furthermore, both of these bottlenecks are resolved: with empirical centralization, in many instances
20 we reduce the asymptotic gap between upper and lower SD bounds, and in this sense, our SD, and thus
21 excess risk, bounds are more transparent. Thus centralization and localization are complementary, and akin
22 to first and second moment normalizations. We intend to fully explore this connection in future work.

23 **Monte-Carlo approach** As mentioned on line 183, a consequence of Theorem 5 is that $n = 1$ trials are
24 sufficient to asymptotically match the rate for the SD shown by Bousquet (2002) (reported in Theorem 3),
25 whereas a use of McDiarmid's inequality (as in the state of the art w.r.t. Monte-Carlo Rademacher Averages
26 (Bartlett and Mendelson 2002)) would require high $n$ for low wimpy variance. Matching Bousquet's rate only
27 occurs with centralization: non-centralized Rademacher averages have inferior Monte-Carlo concentration
28 properties. We have additional experimental results showing a comparison between our results and the one
29 using McDiarmid's inequality, which we did not include due space limitations but we will include in the
30 updated supplementary materials.

31 **RA without absolute values** The identity

$$\mathbb{E}_\sigma \sup_f \frac{1}{m} \sum_{i=1}^m \sigma_i f(x_i) = \mathbb{E}_\sigma \sup_f \frac{1}{m} \sum_{i=1}^m \sigma_i \left(f(x_i) - \hat{\mathbb{E}}_x[f]\right)$$

32 does not hold in general. It does hold by linearity when all functions in the family have equal expectation, but
33 consider the counterexample family $\{x \mapsto -1, x \mapsto 1\}$: any centralized Rademacher average (with or without
34 absolute values) is then 0, whereas without centralization, the value is $\Theta(1/\sqrt{m})$ (see Eq. 5). Other less
35 trivial counterexamples are possible but too convoluted to be explained in the limited space of this response.

36 Thus even the notion of Rademacher averages without absolute values benefits from centralization. Further-
37 more because the 2-sided symmetrization inequality (with centralization and absolute values, eq. 4) is factor-4
38 sharp, and by applying the non-absolute symmetrization inequality once in both directions, we recover a
39 2-sided guarantee, we believe that any gains obtained by removing the absolute value are marginal.

40 We also have evidence that, when using Rademacher averages for non-statistical-learning-theory tasks, such
41 as approximation algorithms for data analytics tasks (e.g., references [17,18,20,21]), the absolute value is
42 important, and should be considered as early as possible.

43 We omitted a discussion of non-absolute Rademacher averages for clarity of presentation, but we could include
44 this information if it is deemed sufficiently important.

45 **Other comments** As requested by Reviewer 4, we will give additional details about applications; in particular,
46 we will describe how our linear family bounds can be applied to get sharper bounds for selecting the optimal
47 expert in a batch-learning panel-of-experts setting (and thus reduced regret in online settings). Reviewer 2
48 asks how our methods extend to model selection, and we note that, for example, this panel-of-experts setting
49 can be extended to *structural risk minimization* if the experts are organized into concentric groups. We will
50 use Big-Oh notation whenever possible, and clarify the statement of Thm. 4. We will fix the typos pointed
51 out by the Reviewers and do a deep editing pass for any other spelling, grammar, or syntax issues.