

1 We thank all the reviewers for their helpful corrections and feedback. We will incorporate them into the revised version.

2 **R1:** We thank R1 for pointing out the missing eigenvalue notation. With regards to comments in weakness section,

3 1. In the fully agnostic setting (where the model is completely misspecified) we don't believe our results apply
4 because the identity of the underlying shared representation becomes ambiguous. But for small degrees of model
5 misspecification, we believe our results can be extended there at the price of additional error terms in our guarantees.

6 2. Our current results assume homogeneous sample size across different tasks for a clean presentation. We believe
7 our techniques/results can be easily extended to non-homogeneous sample sizes across different tasks (with additional
8 notational modifications in our results).

9 3. The “classifier head” of Example 3 (in Appendix A) is actually a nonlinear function. We agree that further
10 investigating the task diversity definition for different examples is an important direction for future research.

11 4. We believe the neural network example will generalize to the logistic loss (by combining with the contents in Section
12 4.1). We opted for our current choice of different losses/functions to exhibit the utility of our general framework.

13 5. We believe that other bounds on the Gaussian/Rademacher complexity can directly be applied in our framework
14 with little modification (such as [1, Theorem 5] for example). However, our framework does not directly accommodate
15 margin bounds—such as in [2]—and extending our results to include these is an interesting future direction for research.

16 6: We agree that investigating the experimental implications of our bounds is an important direction for future work.

17 **R2:** We thank R2 for pointing out the typos, and suggestions for clarification. With regards to the typos in the
18 “Additional Feedback” Section. For 1: $\hat{\mathbf{f}}$ refers to an ERM solution in the task variables of the training risk in Eq.
19 2—the pair $(\hat{\mathbf{f}}, \hat{\mathbf{h}})$ refers an entire ERM solution in Eq.2; we will define/clarify this. For 2: The union should be written
20 $\mathbf{h} \in C_{\mathcal{H}_X}$ (referring to the \mathcal{H} covering with respect to inputs \mathbf{X}). For 3: Yes, both terms should be primed (and the
21 statement should be more clearly written as, “given this \mathbf{h}' , $\exists \mathbf{f}' \in C_{\mathcal{F}_{\mathbf{h}'(\mathbf{X})}}^{\otimes t}$ that is ϵ_2 -close to \mathbf{f} with respect to inputs
22 $\mathbf{h}'(\mathbf{X})$. By construction, $\mathbf{h}' \in C_{\mathcal{H}_X}$ and $\mathbf{f}' \in C_{\mathcal{F}^{\otimes t}(\mathcal{H})}$.”). We will correct/clarify the notation here.

23 With regards to Gaussian complexity vs Rademacher complexity, we chose to include only Gaussian complexities in the
24 main paper to simplify the presentation, since the chain rule is most naturally stated in terms of Gaussian complexities.
25 Since Gaussian/Rademacher complexities are equivalent up to logarithmic factors [3, p.97] we could also rewrite all our
26 results in terms of Rademacher complexities at the cost of only logarithmic factors.

27 With regards to the comments in the weakness section, we believe our notion of task diversity can be understood in a nat-
28 ural way, and will provide further discussion and intuition to interpret it—see the following. Our framework/arguments
29 (which hold for general \mathcal{F} , \mathcal{H} and ℓ) do require abstraction, but we also believe this abstraction is a strength that allows
30 our guarantees to be applied to a wide class of problems.

31 Definition 1 and Definition 2 seek to define two notions of distance between two representations \mathbf{h}, \mathbf{h}' . In our
32 framework, information about the representations is only observed through the composite functions $f \circ \mathbf{h}$. For any
33 direction/component in \mathbf{h} that is not seen by a corresponding task f , that component of the representation \mathbf{h} cannot
34 be distinguished from a corresponding one in a spurious \mathbf{h}' . When this component is needed to predict on a new task
35 f_0 which lies along that direction, transfer learning will not be possible. Therefore, Definition 1 defines a notion of
36 representation distance in terms of information channeled through the training tasks while Definition 2 defines it in
37 terms of an arbitrary new test task. Task diversity essentially encodes the ratio of these two quantities (i.e. how well
38 the training tasks can observe relevant parts of the representations useful for the new task). In the case where the \mathcal{F}
39 contains underlying linear task functions $\alpha_j^* \in \mathbb{R}^r$ (as in our examples in Section 4), our task diversity definition
40 reduces to ensuring these task vectors span the entire r -dimensional space containing the output of the representation
41 $\mathbf{h}(\cdot) \in \mathbb{R}^r$. This is quantitatively captured by the conditioning parameter $\tilde{\nu} = \sigma_r(\mathbf{A})$ for $\mathbf{A} = (\alpha_1^*, \dots, \alpha_t^*)^\top \in \mathbb{R}^{t \times r}$
42 which represents how correlated these vectors are in \mathbb{R}^r . Appendix A gives a further task diversity example when \mathcal{F}
43 contains nonparametric functions. We will provide further explanation of these definitions and their relationship to each
44 specific example in the final version.

45 **R3:** We thank R3 for their comments and agree studying transfer learning in frameworks other than those using
46 Gaussian complexities (i.e. with more refined data-dependent bounds as R1 mentions), is an interesting future direction.

47 [1] Golowich, Noah, Alexander Rakhlin, and Ohad Shamir. “Size-independent sample complexity of neural networks.”
48 Conference On Learning Theory. 2018.

49 [2] Wei, Colin, and Tengyu Ma. “Improved sample complexities for deep neural networks and robust classification via
50 an all-layer margin.” International Conference on Learning Representations. 2019.

51 [3] Ledoux, Michel, and Talagrand, Michel. Probability in Banach Spaces.