

1 We individually respond to the questions and concerns raised by reviewers about our work below:

2 **Reviewer 1:**

- 3 • **"The framework is natural but straight-forward"**: We are glad that R1 thinks that our framework is natural.
4 However, inspite of being simple, such a unified view has not been applied to explain popular representation learning
5 approaches, especially recent ones such as masked self-supervision and VAE. We believe that formulating these
6 approaches theoretically and analysing their sample complexity bounds is an important problem in the domain of
7 representation learning and can lead to insights. We view the simplicity of the problem formulation to be a strength
8 rather than a weakness.
- 9 • Additionally, we present two concrete applications of the framework, with empirical support thereby showing the
10 validity of our framework. This shows that our general framework is effective and doesn't lead to vacuous bounds.

11 **Reviewer 2:**

- 12 • **"Perhaps to have more focus on cases when this will not work"**: Our theorems and results do provide implications
13 for cases when the auxiliary self-supervised task may not help the prediction task. For example consider Theorem 1,
14 if $\mathcal{H}_{\mathcal{D}_X, L_r}(\epsilon_0)$ is not significantly smaller than \mathcal{H} , then using unlabeled data will not reduce the sample size of the
15 labeled data much compared to only using the labeled data for prediction. Further, for the two concrete examples
16 that we present in Section 5, the sample complexity bounds also indicate when the unlabeled data will not be very
17 useful. Another possible reason that these representation learning approaches can fail is that the optimization does
18 not provide a good solution. However, accounting for this is out of scope of this work. We will incorporate the
19 suggestion and add more discussion about the related works pointed out.
- 20 • **Clarification about Equation 15**: This equation means that if we fix a pair (h, g) with $L_r(h, g, D_x) \geq \epsilon_0$, then
21 we have $P(L_r(h, g, U) = 0) \leq \delta/2|H||G|$. Since there are at most $|H||G|$ such (h, g) pairs, by the union bound
22 we have $P(\exists(h, g) \text{ s.t. } L_r(h, g, D_x) \geq \epsilon_0 \text{ and } L_r(h, g, U) = 0) \leq \delta/2$. Then with probability at least $(1 - \delta/2)$,
23 there exists no (h, g) such that $L_r(h, g, D_x) \geq \epsilon_0$ and $L_r(h, g, U) = 0$. This means that only those (h, g) with
24 $L_r(h, g, D_x) \leq \epsilon_0$ will have $L_r(h, g, U) = 0$.

25 **Reviewer 3:**

- 26 • **"Interaction between the two learners"**: Our sample complexity bounds do quantify the interaction between the
27 two learners (from the hypothesis classes \mathcal{H} and \mathcal{G}) by introducing the notion of $\mathcal{H}_{\mathcal{D}_X, L_r}$. This captures the effect
28 that the representation learner over $\mathcal{G} \circ \mathcal{H}$ has on the prediction learner over $\mathcal{F} \circ \mathcal{H}$.
- 29 • **"Covering numbers are not very effective in practice"**: While we acknowledge the evidence that naively applying
30 uniform convergence bounds may not result in good generalization/sample bounds for deep learning, these existing
31 studies may not apply to the setting we consider in our paper. To the best of our knowledge, this evidence is specific
32 to supervised learning without the auxiliary representation learning tasks, while our setting is with the auxiliary tasks.
33 In particular, the mentioned paper "Uniform convergence may be unable to explain generalization in deep learning"
34 does not apply directly to VAE or any other representation learning approach studied in our work.
35 Furthermore, our experimental results do correspond with our sample bounds in Section 5 which are based on uniform
36 convergence. This is in contrast to the existing studies on supervised deep learning without auxiliary representation
37 learning tasks. These results suggest considering a shifted view: "Uniform Convergence strikes back and can explain
38 the generalization behavior of deep learning *with* auxiliary representation learning tasks". Why so? Without
39 auxiliary representation learning tasks, it is generally believed that the optimization has an implicit regularization on
40 the training, and hence uniform convergence fails to explain this. However with the auxiliary tasks, we conjecture
41 that functional regularization restricts the learning dynamics to a smaller subset of hypotheses, on which the implicit
42 regularization of the optimization is no longer significant, and thus the generalization can be explained by uniform
43 convergence. This is an interesting open question and we leave it as future work.

44 **Reviewer 4:**

- 45 • **"Experiments on real world data would have helped"**: We have presented some experiments on real data in
46 the appendix. Due to space restrictions, we only focus on experiments for the two concrete instantiations of our
47 framework, which we believe can give fine-grained empirical evidence (e.g., how the sample bounds depend on r ,
48 etc.) for our analysis. We will move some experiments on real data to the main body in a future version of our paper.
- 49 • **Line 135**: The loss $l_c(f(h(x)), y)$ denotes the loss function for the prediction task, and has been defined on line 90.
- 50 • **Line 179**: By "standard analysis", we refer to the standard statistical learning theory argument for uniform conver-
51 gence over a finite hypothesis class. We will make this explicitly clear in in a future version of our paper.
- 52 • **Related work**: Thanks for the suggestions! We will add them to the future version.