1 We thank all the reviewers for their insightful comments and suggestions that will all be incorporated in the final version.
2 We start our response with Reviewer #3, since some of the comments made there are inaccurate.

3 **Response to Reviewer #3**: The major critique is that our paper rediscovers prior ideas. However, the two papers of
4 Ben-David and others pointed out by the reviewer are essentially different in the setups and results, which is also why
5 those papers were not cited in more related works we cite (e.g., [11,17,31]) either. Below we expand these differences:

6 1. The setting. They assumed a *multitask learning* setting where $n$ tasks/datasets symmetrically contribute to the loss
7    function, and the learner has access to *datasets for all tasks* (i.e., theorem 3.5 assumes $|S_i| \geq m$ for all $i$). In contrast,
8    our framework assumes an asymmetric *supervision* setting where the learner only has access to samples of $(X, O)$
9    but *no dataset for the target label $Y$*. Moreover, They assumed that all tasks are *binary classification* while in our
10    framework $\mathcal{Y}$ and $\mathcal{O}$ can be *multiclass* and different. Also, they define the relatedness of tasks by deterministic
11    *functions $X \rightarrow X$* while in we define it as *conditional distributions $\mathbb{P}(O|Y)$*.
12 2. The analysis and results. In their work, since all tasks are associated with a dataset of size $\geq m$, the generalization is
13    given for-free, as long as the generalized VC-dimension is finite. However, in our paper, learnability is nontrivial
14    (e.g., think about a very noisy dataset). So, our paper makes additional effort to characterize the transition class and
15    the learner's prior knowledge about it, which results in conditions [C2] and [C3].

16 The other critique is we make strong assumptions. In fact, our assumptions are *common* and often *weaker* than literature.

17 1. Realizability assumption. This is commonly assumed in the related works (see line 34) and is approximately
18    satisfied by many practical models (such as DNNs that achieve classification error close to 0 on many large, practical
19    problems). Due to the stochastic generating process of the indirect signal, it would be technically hard to remove
20    this assumption in a general setting. We therefore think that it is a good starting point for developing the theory.
21 2. Learnability conditions. First, the conditions [C2] and [C3] actually generalize and relax many standard learning
22    conditions in the literature (see Example 5.6 & 5.7). Moreover, the second part of our theorem 5.2 shows that the
23    conditions are not only sufficient, but also close to be necessary. In other words, our framework provides more
24    practical (rather than unrealistic) ways to determine learnability and/or find learnable supervision signals.

25 **Response to Reviewer #1**: Here are our responses to the concerns:

26 1. Use of VC-style argument. We choose to use VC-style argument to bound the Rademacher complexity because it is
27    a *safe* way to show *learnability*. It is safe in the sense that most of the practical machine learning models, including
28    DNNs, have a finite VC-dimension. Also, it is possible to adopt to a PAC-bayesian argument by defining induced
29    prior and posterior via the transitions. However, to show learnability in this way, we would need to further bound the
30    KL-divergence by constructing priors for specific models. It would be an interesting direction for future research.
31 2. Practical implications. We will address this issue more clearly in the final version. Briefly: To make full use of the
32    indirect observations, our framework provides a practical and unfied way to understand the supervision power of the
33    signal when only weaker form of the prior knowledge is given (encoded by *separation*). Our framework generalizes
34    previous results, and hence suggests new learning scenarios. Also, in a non-learnable case, section 5.3 also guides to
35    find complementary signals to make the problem learnable.

36 **Response to Reviewer #2**: We thank the reviewer for the useful feedback. For the weakness and additional feedback:

37 1. To derive a matching lower bound, we will need two things: (1) For the case when $\mathcal{T} = \{I\}$, it is true that faster
38    rate can be achieved. This is partly because the task $X \rightarrow O$ is agnostic in default. So if using the current method,
39    (at least) we will need to define what is a "realizable" $O$. (2) A new measure similar to separation that links the
40    annotation and classification risk. We will add more discussion on this issue and/or the limitations fo current work.
41 2. We will add an illustrating example for this phenomenon 199-202. For linear case, the intuition is that different
42    regions will have different patterns of the distribution of $O$. We will also add a visualization to make it clear.
43 3. For the additional feedback: **(1,5)** There are typos. Will be corrected in the complete version. **(2)** We will adopt
44    the suggested way of definition and call it "$\mathcal{T}$-learnable". **(3)** $b$ is the upper bound of the annotation loss $\ell_{\mathcal{O}}$, which
45    may not be the 0-1 loss. We will add a reminder here. **(4)** We will add a definition for the Natarajan dimension.
46    **(6)** We can restate the definition as: $\mathcal{D}_i(x)$ is the set of all the induced distributions $\mathbb{P}_T(O|X = x, Y = y_i)$. **(7)**
47    $\dim$ is used to indicate the singularity of $T$, since one may concern if a singular $T$ will lead to information loss. **(8)**
48    Our framework can be extended to product case and the product label will contain more information and provably
49    have better supervision power. However, to formulate this idea in general, we also need to model and characterize
50    the *correlations* between different annotations (which could be expensive). It will be an interesting future research
51    direction. **(9)** Your intuition is correct. We would describe it like this: [C2]: the optimal classifier of $Y$ can induce
52    an optimal predictor of $O$. [C3]: the suboptimal classifier of $Y$ will induce a suboptimal predictor of $O$.