

1 We thank the reviewers for their insightful feedback. We are pleased that the reviewers found our paper well-written
2 [R1,R2,R4], thorough [R1,R3], and recognized the strength of our theoretical results [R3,R4]. We are especially
3 encouraged by the comments of R3 and R4 concerning the possible impact of our work beyond the information
4 bottleneck (IB). We first address two recurrent questions and proceed by addressing reviewer-specific comments.

5 **G.1. Clarifying the definition of $\text{Dec}(X, y)$ [R1,R2,R4].** We simplified the notation and added an example to better
6 illustrate this concept. We summarize these changes here. Suppose we are classifying cats versus dogs, then the y (say,
7 “cat”) decomposition of images X , corresponds to taking all “cats” and labeling them with every possible binary labels.
8 Formally, let X_y be a r.v. s.t. $P_{X_y} = P_{X|y}$ (the “cat” r.v.). Then $\text{Dec}(X, y) := \{N|\exists t' : \mathcal{X} \rightarrow \mathcal{Y} \text{ s.t. } N = t'(X_y)\}$. We
9 emphasize that: (i) N is a “labeling” of X_y so it takes value in $\mathcal{Y} (\{\text{“cat”}, \text{“dog”}\})$ so we can predict it using a binary
10 classifier in \mathcal{V} ; (ii) the formal definition does not assume that the underlying labeling of Y is deterministic.

11 **G.2. Practicality of our Decodable IB (DIB) [R1,R2].** We extended the discussion about practical optimization
12 (summary of Appx. D6, D7) and estimation (summary of Appx. D4) in the main paper. We emphasize that: (i) the
13 min-max problem can be well optimized using joint gradient descent ascent (c.f. Fig. 10 and L. 907), it does not require
14 unrolling steps; (ii) we get good performance by only predicting a few (4) “labels” $N \in \text{Dec}(X, y)$ (c.f. Appx. D6). As
15 a result, training a network with DIB is approximately as computationally efficient as training a standard network or the
16 variational IB — and much more efficient than IB. We will **release our code** to help practitioners use DIB.

17 **R.1.** We would like to thank you for the helpful review and useful comments which we addressed in the revised
18 manuscript. We are pleased that you found the idea interesting, well-motivated, and the analysis thorough.

- 19 • “Can we prove such representation exists?”: **Yes**, they always exist because they are defined as maximizers/minimizers
20 of some quantity in a finite set. This is a very good point, which we now discuss in the paper.
- 21 • “[...] I found part 3.2 unclear [...]”]: Please refer to G.1. for $\text{Dec}(X, y)$ and N . We also added an algorithm on the
22 optimization of $I_{\mathcal{V}}[Z \rightarrow \text{Dec}(X, Y)]$. For cat-dog classification: (i) assign each cat with a binary label; (ii) optimize
23 the encoder such that classifiers in \mathcal{V} cannot predict these labels, i.e., distinguish between cats; (iii) repeat for dogs.
- 24 • “[...] Does the method induce additional complexity [...]?”]: **No**, please refer to G.2.
- 25 • “Consistency between [...] proposition 1. [...] Fig. 3 [...]”]: We would like to clarify that Prop. 1 only concerns the
26 existence of an optimal predictor. The generalization of any ERMs requires \mathcal{V} -minimality (Theorem 1). The best loss
27 in Fig. 3a occurs at approximate \mathcal{V} -minimality (Fig. 3b) suggesting that ensuring generalization is more valuable
28 than improving the best possible performance. We thus do not believe that Fig. 3 contradicts the theory.
- 29 • “Usefulness of the method: [...] traditional IB [...]”]: Due to space constraints, we decided to highlight results
30 that validate the theory, and thank the reviewer for acknowledging the importance of such results. We nevertheless
31 emphasize that many of the results asked by the reviewer are in the appendices. Table 3 (Appx. D10) evaluates DIB
32 as a regularizer on CIFAR10MNIST (CIFAR10 with overlaid MNIST images to evaluate robustness to “nuisances”) and
33 MNIST (to replicate the results from VIB’s paper). We see that DIB outperforms VIB on both tasks, i.e., it gives
34 rise to a better classifier and is more robust to nuisances. We added a small discussion about this in the main text.

35 **R.2.** We thank you for the kind and helpful review, we are pleased that you found the paper well-written and the
36 approach elegant. We have put significant effort into incorporating both your suggestions into the revised manuscript.

- 37 • “[...] the main text could benefit [...] practicality of DIB compared to IB [...]”]: Please refer to G.2.
- 38 • “[...] more clear, maybe through an example, is in explaining equation 4. [...]”]: Please refer to G.1.

39 **R.3.** We thank the reviewer for reading the paper thoroughly and for providing a kind and detailed assessment. We are
40 encouraged that you found the idea novel, the theoretical results strong, and the empirical evaluation thorough.

- 41 • “Although Theorem 1 is a strong result, it assumes [...]”]: We believe that our assumptions **do not significantly limit**
42 **the applicability of our theory**: (i) Piece-wise universality: we now **relaxed this assumption** so that the proof
43 holds for any model using a softmax layer with biases on that layer; (ii) Finite sample spaces: although the continuous
44 setting is theoretically interesting, it (arguably) is less important in practice as everything must be discretized on a
45 computer; (iii) Deterministic labeling: examples are usually seen once per dataset, i.e., no two same images are given
46 different labels. In such cases the labeling is deterministic. We acknowledge that this is not necessarily true in the
47 real world and hope to extend our proof in future work. We emphasized this assumption in the revised paper.

48 **R.4.** We thank you for your helpful review, and are pleased you think the approach could constitute a significant
49 advancement for IB and be broadly relevant for methods based on information theory (e.g. variants of VAE and GAN).

- 50 • “The choice of the appropriate predictive family \mathcal{V} [...]”]: There is a trade-off. A larger \mathcal{V} means a smaller best
51 achievable loss, but a more complex \mathcal{V} makes the estimation of \mathcal{L}_{DIB} provably harder. This is a very good point,
52 which we now discuss in the revised paper.
- 53 • “[...] the decomposition of the input X [...] might need a revision”]: Please refer to G.1.
- 54 • “What is [...] odd/even CIFAR100?”: Classifying the parity of the class index. This “CIFAR2” is easier to visualize.