1 We thank the reviewers for their substantive and constructive feedback, and appreciate their assessment of our explanation

2 framework as "simple yet effective" **(R1)** with convincing experiments and clear presentation. We particularly appreciate

3 **(R4)**'s comment that our "detailed and thoughtful approach presents a great example of high-quality reserch." The main

4 reviewer concerns were narrative-related; we have addressed these concerns (and other specific suggestions) in revision.

5 **1: Technical misconceptions in reviews.**

6 • **(R1)** (Weakness 3) Section 4 illustrates and tests the performance of our method in a controlled setting; here we

7 define the data distribution to be Gaussian, but outside of this section we do *not* make any assumptions on the data

8 distribution. Also, our objective does *not* attempt to minimize $I(\beta; \widehat{Y})$ (but see Appendix A for discussion of potential

9 objective variants).

10 • **(R1)** (Weakness 4) We do *not* require the training data originally used to train the classifier, which we agree would be

11 problematic. Our explanations are based on a data distribution, but this need not be the training distribution.

12 • **(R1)** **(R2)** The validity of our DAG stems from the independence of $(\alpha, \beta)$, which we impose with an isotropic normal

13 prior. (We have clarified this point in Sections 3.3, 4, and 5). Although MI is indeed a correlative metric in general,

14 the independence of $(\alpha, \beta)$ allows us to show that in our framework it quantifies information flow, a well-founded

15 node-based metric for causal influence (see refs 7, 55). This argument is indeed very simple (as pointed out by

16 **(R1)** **(R2)**), but it depends crucially on our modeling decisions. Note that in the generative modeling literature the

17 assumption of latent factor independence is common. One might alternately allow dependencies between the latent

18 factors, jointly learning their causal structure with the generative network. However, without labeled side-information

19 that would give these features semantic meaning, in our view our independence-by-construction approach is more

20 useful to generate parsimonious explanations.

21 **2: Changes to storyline.** The most consistent reviewer concerns were narrative-related: (1) **(R1)** **(R2)** **(R4)** com-

22 mented that the role of causality in our framework was overstated or unnecessary (since by enforcing the $(\alpha, \beta)$ to be

23 independent our metric for causal influence simply reduces to MI), and (2) **(R1)** **(R2)** **(R3)** **(R4)** commented that our

24 narrative did not engage with the disentangled representation literature. We thank the reviewers for these suggestions

25 and have addressed them with careful changes to our storyline and a major addition to our related work section:

26 • As suggested by **(R1)** **(R4)**, we have adjusted our storyline to change the perspective from which we approach

27 causality, instead highlighting the aspects of the model that allow the simple MI metric to be interpreted causally

28 **(R2)**. We also clarified that the prior enforcing independence of $(\alpha, \beta)$ leads to the validity of our DAG **(R1)**. We

29 agree with **(R4)** (and related work) that MI on its own can serve as a valid theoretical justification for explanation,

30 and have reworked our discussions of MI-based methods to remove the impression of downplaying, emphasizing

31 instead that our framework provides the *complementary* benefits of causal and information-theoretic interpretations.

32 • We believe that the disentanglement framing suggested by **(R1)** **(R2)** **(R3)** **(R4)** is exciting and could open new

33 research directions. Our method can indeed be thought of as an information-based disentanglement procedure. Unlike

34 techniques such as InfoGAN, our method (1) uses a classifier as side information; (2) separates classifier-relevant

35 from classifier-irrelevant features, making the framework suitable for explanation; and (3) allows the MI metric to be

36 interpreted as a measure of causal influence of disentangled features on the classifier output. We have adjusted the

37 storyline throughout the paper to add this perspective, and have added a paragraph in the background to discuss the

38 relation to specific disentanglement techniques (including all suggestions from **(R1)** **(R2)** **(R3)** **(R4)**) in more detail.

39 **3: Additional experiments and discussion.**

40 • **(R4)** We agree that computational cost is a concern. We have added a sentence in the main text mentioning this

41 cost, and have expanded discussion in the broader impacts section of the potential drawbacks of using complex,

42 uninterpretable models for explanation.

43 • **(R3)** **(R4)** The question of how mismatched explainer model capacity affects results is indeed both interesting and

44 important. We address this in two ways. First, we are currently performing an experiment in which latent space

45 dimension and VAE architecture complexity are swept for a fixed classifier complexity. We will add these results

46 (presented both qualitatively with sample explanations, and quantitatively as in Figure 5(a-b) and Supplement Figure

47 11) along with a brief discussion in the main text. Second, we can use results from Feder & Merhav (*IEEE Trans.*

48 *Info. Theory*, 1994) to bound the capacity mismatch of our explainer (i.e., explainer error in predicting classifier

49 outputs) with the $I(\alpha; \widehat{Y})$ part of our objective. In practice, this result means that a sufficiently large value of $I(\alpha; \widehat{Y})$

50 serves as a certificate that the explainer complexity is sufficient to explain the classifier. We have added an appendix

51 containing details and implications of this analysis and a brief discussion to the conclusion of the main paper.