1 We thank the reviewers for their insightful and valuable feedback and for their *unanimous support* of the paper.

2 We are encouraged that they found our formulation to be **"novel/new"** (**R1**,**R2**,**R3**) **"interesting"** (**R2**,**R3**,**R4**) and

3 with a **"strong direction"** (**R4**); The experimental results to be **"effective"** (**R1**,**R2**,**R3**,**R4**), **"extensive"**,**"thorough"**

4 (**R2**,**R4**), **"detailed"** (**R2**,**R3**) and **"strong"** (**R4**); And AO-CLEVr dataset **"beneficial"**,**"good"** (**R2**,**R3**).

5 **R4: DAG is most suitable to disentangled attributes, some attribute manifest differently depending on object.** Extend-

6 ing to dependent pairs is a very important next direction. Unfortunately, it is not clear that zero-shot can work well

7 with strongly entangled pairs because every case could be special. Three insights worth mentioning: (1) Even the

8 fully disentangled case is still very challenging. (2) The factored DAG can be used as a strong zero-shot prior for

9 few-shot learning, thus benefiting future work on dependent pairs. (3) Using a "closed" settings may capture some of

10 the dependency by eliminating "over-generalization" (e.g. by disallowing yellow-wine label).

11 **R4: Clarify details of mapping the causal graph to Fig 1c (1) The role of $g_A^{-1}, g_O^{-1}$, do they add assumptions? context**

12 **to causal DAG?** $g_A^{-1}$ and $g_O^{-1}$ are used to estimate the latent $\phi_a$ and $\phi_o$ of an image instance. They reflect as-

13 sumptions about the noise level in the data-generation process (Suppl L508-513), i.e. that the mapping from the

14 core-features ($\phi_a$ and $\phi_o$) to the image ($\mathbf{x}$) is not too noisy and the latent vector can be recovered from the image.

15 **(2) The new nodes $\hat{\phi}_a, \hat{\phi}_o$ satisfy the independence constraints by construction. Explain consistency.** We respectfully

16 point out that since $\hat{\phi}_a, \hat{\phi}_o$ are children of $\mathbf{x}$, they *do not* satisfy the independence constraints of Eq. 6. Minimiz-

17 ing $L_{indep}$ encourages the property $p(\hat{\phi}_o|do(o)) \approx p(\hat{\phi}_o|do(a, o))$ (L550). Only then the independence relations of Eq.

18 (6) apply to $\hat{\phi}_a, \hat{\phi}_o$. It also minimizes the PIDA metric of (Suter 2019). **(3) Any assumptions fail for MLPs?** No.

19 **R4: Where does the causal interpretation manifests? (1) Difference from standard embedding?** As the reviewer points

20 out, one difference is in the independence loss; another is the use of two separate embedding terms tied to the in-

21 dependence loss; both motivated by the causal graph. We deliberately proposed a model close to baselines (L186)

22 to surgically demonstrate the strength of the proposed approach. **(2) Why is $\lambda_{invert}$ essential?** Since there exist no

23 ground truth values for neither $\phi_a$ nor $h_a$, minimizing $||\hat{\phi}_a - h_a||^2$ may reach trivial solutions (same for $\phi_o, h_o$).

24 $\lambda_{invert}$ guides the optimization and avoids trivial solutions. It does not contradicts assumptions on the causal process.

25 **R1: How are the means $h_a, h_o, g(h_a, h_o)$ updated?** Instead of learning explicit values for the means, we learn MLPs

26 that output the means (using gradient updates L138,143,178). For example, an MLP ($h_A$) maps the (one-

27 hot) representation of "leather" to $h_{leather}$ and an MLP ($g$), maps ($h_{leather}, h_{sandal}$) to $g(h_{leather}, h_{sandal})$.

28 **R1: Does interventional inference means matching prototypes?** Partially yes: Inference that follows *the approxima-*

29 *tions we took* (Supp A., e.g. Gaussian and $0^{th}$ order Taylor) may be viewed as matching prototypes. In the general

30 case, there may be better ways to estimate the likelihood of $p(\mathbf{x}|a, o)$ and the factors $p(\phi_a|a), p(\phi_o|o), p(x|g(\phi_a, \phi_o))$.

31 **R1: Disentanglement is achieved by independence loss rather than intervention.** The independence loss allows to

32 learn a model that is robust to interventions. Minimizing $L_{indep}$ encourages $p(\phi_o|do(o)) \approx p(\phi_o|do(a, o))$ (L550).

33 **R2: Independence loss encourages the performance on the unseen data but drops on the seen data.** This is a known

34 and important trade-off (Rothenhäusler 2018): The independence loss discourages certain types of correlations, hence

35 models do not benefit of them when the test and train distributions are identical. However, the loss is constructed

36 in such a way that these are exactly the correlations that fail to hold once the test distribution changes (zero-shot).

37 Ignoring these correlations improves performance on unseen data. We will refer to (Rothenhäusler 2018) and discuss.

38 **R2: Failure analysis.** Following this request, we analyzed samples of unseen pairs of Zappos in the open-world setup.

39 We compared *Causal* with LE*, which is the strongest no-prior baseline. LE* confuses unseen pairs for seen pairs

40 at a rate of $3.7{:}1$, while *Causal* errors are more balanced $1.2{:}1$. One interesting failure case of *Causal*, is that it

41 over-commits for predicting the pair "Leather-Slippers", which was unseen during training. In the final version we

42 will provide more qualitative and quantitative details about Zappos and AO-CLEVr .

43 **R3:MIT dataset: Consider top-k labels**: Following this suggestion, we conducted a new experiment to evaluate both

44 top-1 and top-2 accuracy. Raters were asked to select the best and 2nd-best attributes that describe an image, among

45 attributes relevant for that object. The top-1 accuracy was 32%, consistent with previous experiment. The top-2

46 accuracy was 47%, only slightly higher than adding a random label on top of top-1 label (yielding 42%). To verify

47 that raters were attentive, we also injected 30 "sanity" questions that had two "easy" attributes, yielding top-2=100%.

48 **R3: MI instead of HSIC:** HSIC advantage is that it is non-parametric, unlike MI, and does not requires training

49 an additional network for variational approximation. **Embed. size**; We will report results w.r.t. embedding size.

50 **Efficacy of HSIC:** See $\lambda_{indep}$=0 at Table S.2. **Results w/o alternate training.** Alternate training lowers the SEM (L713).

51 Means are comparable (68.7 vs 67.7).

52 **R3: Revise method for smoother reading. R4 Some bits are confusing.** We will restructure the paper based on your

53 feedback: (1) Shorten the "overview" section (2) Discuss how independence loss allows $\hat{\phi}_a, \hat{\phi}_o$ to recover the proper-

54 ties of $\phi_a, \phi_o$, and its relation to PIDA. (3) Update the final version based on the rebuttal.

55 **R1**,**R2**,**R3**,**R4**: We will address all minor comments, and clarify the broader impact.