

1 We thank the reviewers for their detailed feedback. We're glad that the reviewers thought our general identification
 2 theorem is sound and interesting (R1, R2, R3), found our work is creative (R2), and see VDE as an elegant solution to
 3 satisfy the independence conditions required in theorem 1 (R3).

4 To address clarity concerns, we have adopted the reorganization suggestions of section 2 from R2. This should move the
 5 "identification comparison with other control function methods" closer to theorem 1 which clarifies one of R1's main
 6 concerns by connecting our work to Guo and Small 2016. This would also surface the error analysis of GCFN from
 7 section 2.3 and appendix A.7 that R1 requests and provide the more intuitive structure of section 2 that R3 requests.

8 **[R1 + R3 : limited experimental evidence | R3, toy experiments]** We do not know of any established benchmarks
 9 for causal effect estimation with IVs. We briefly discuss our evaluation here. Section 3.1 shows GCFN can handle
 10 assumptions that existing methods like CFN, 2SLS, and DeepIV cannot. Then, as demonstrated below and in section
 11 3.3 in the paper, GCFN performs on par or better than DeepIV and DeepGMM on *high-dimensional* simulated data
 12 from each paper [1, 2] respectively. Thus, GCFN is competitive on data that satisfies the assumptions required by
 13 DeepIV and DeepGMM *while* being applicable to data generated with non-additive outcome processes [1, 2].

14 **[Additional evaluation of GCFN]** Due to reviewers' concern about toy experiments, we present further evaluation
 15 of GCFN on simulated data with a high-dimensional IV. We use the data generating process given in DeepGMM [2], a
 16 recent state of art method. We ran GCFN with 10 different random seeds and report results for $\kappa = 0.3$, chosen based
 17 on mean outcome MSE. We report results for DeepGMM and DeepIV as reported in [2]. GCFN performs competitively
 18 with an effect MSE of 0.077 ± 0.022 compared to DeepGMM's 0.07 ± 0.02 and DeepIV's 0.11 ± 0.00 . Effect
 19 MSE for $\kappa \in \{0.2, 0.4\}$ were similarly within standard errors of DeepGMM's performance and better than DeepIV [1].

20 **[R1 + R3, hard to gauge assumptions | realism of assumptions]** Solving VDE exactly guarantees reconstruction
 21 (A1) and marginal independence $\hat{z} \perp \epsilon$. A strong IV (A2), intuitively, is one that is able to set treatment to any value
 22 given any fixed confounder value. Joint independence requires assumptions on the treatment process like additivity.
 23 Additivity is common in the economics literature and strong IV is an assumption that can be reasoned about using
 24 domain expertise: for e.g. can college proximity influence a student's decision to go to college regardless of skill?

25 **[R1, Theorem 1 and traditional identification results]** To clarify, theorem 1 characterizes control functions that
 26 guarantee identification. We compare GCFN's identification for an additive treatment process to traditional control
 27 function identification from Guo and Small (reference 15 in the paper) in the paragraph "New conditions for effect
 28 identification" in section 2.2. Briefly, a strong IV requires more than the exclusion restriction and relevance properties
 29 of IVs. For this added assumption, GCFN drops the additive outcome process and the additional conditions on noise
 30 that CFN requires. We have expanded this discussion in the paper.

31 **[R1, Section 2.2 is opaque, no analysis]** As noted above, we
 32 have clarified our presentation in section 2.2 and A.7. We briefly
 33 summarize GCFN here: 1) GCFN's first stage, VDE, has an ob-
 34 jective where at optimum, control functions meet condition A1 in
 35 theorem 1. Similarly, VDE's decoder leverages treatment process
 36 assumptions where at optimum joint independence (A3) is satisfied
 37 (line 187). 2) GCFN then runs a flexible outcome regression on
 38 the treatment and VDE's control function to estimate effects. As
 39 in section 2.3, we analyse GCFN's effect error in appendix A.7. In
 40 A.7.1, we show how the two components of VDE, reconstruction
 41 error and dependence of \hat{z} on ϵ , influence effect error in data with
 42 additive treatment processes. We also give a general bound in A.7.2.

43 **[R2, $I(\hat{z}, \epsilon) > 0$, reconstruction error, GCFN failure cases]**
 44 We thank the reviewer for this suggestion. We will add discussion
 45 about failures modes of GCFN. We plot the influence of non-zero
 46 information and reconstruction on effect MSE in fig. 1 using models
 47 trained on the MNIST IV data above for different κ .

48 **[R3, Real experiment unconvincing]** The reported effect in Nunn et al. (reference [34] in the paper) is a well-known
 49 result supported by modelling choices informed by domain knowledge. We believe that recovering this effect gives
 50 evidence that GCFN works well even without strong parametric assumptions used by Nunn et al.

51 **[R3, Kernel IV, Multidimensional \mathbf{z}]** We thank the reviewer for the suggestion. We will include a Kernel IV baseline.
 52 If \mathbf{t} is scalar but \mathbf{z} is multidimensional, an additive treatment process with a scalar $f(\mathbf{z})$ would be $\mathbf{t} = g(\epsilon) + f(\mathbf{z})$. In
 53 this case a scalar \hat{z} suffices to capture $f(\mathbf{z})$ and satisfy ignorability.

- 54 [1] Jason Hartford, Greg Lewis, Kevin Leyton-Brown, and Matt Taddy. Deep iv: A flexible approach for counterfactual prediction.
 55 In *Proceedings of the 34th International Conference on Machine Learning-Volume 70*, 2017.
 56 [2] Andrew Bennett, Nathan Kallus, and Tobias Schnabel. Deep generalized method of moments for instrumental variable analysis.
 57 In *Advances in Neural Information Processing Systems*, 2019.

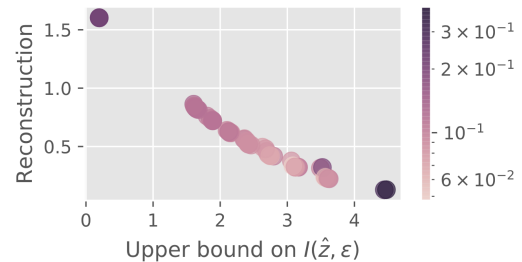


Figure 1: Effect estimation MSE of GCFN in the MNIST IV experiment above plotted for the corresponding VDE's reconstruction loss and upper bound on $I(\hat{z}, \epsilon)$ (as in derivation in A.2); the lighter the color, the better the error. Effect MSE is good and very mildly sensitive to changes in reconstruction or information for values in 0.3 – 0.5 and 2.5 – 3.5 respectively. Outside those ranges, effect MSE is more sensitive and is bad when either quantity is large.