

1 **Reviewers 1 + 3:**

2 | *'Hard to see how TRE was relevant for energy-based modelling ... [the] description of NCE was insufficient'*

3 Thank you for pointing this out; we will remedy this by strengthening the NCE connection from earlier on in the paper.

4 | *'AIS connections could be better'*

5 We agree, and will add further discussion. However, we presently think the connection is more conceptual than technical.

6 | *'MNIST is rather easy. Does it work on harder datasets?'*

7 | *'Would higher dimension problems require considerably more intermediate distributions? If so, how bad is it?'*

8 While too preliminary to be included in the paper, we have some evidence from ongoing work showing that TRE can  
9 work on higher-dimensional image datasets using only modestly more intermediate distributions.

10 | *'Is it easy to say if the limitations in [Stratos & McAllester] apply to your method and, if not, why not?'*

11 We believe the answer is no, since Stratos & McAllester prove strong limitations on “high-confidence lower bounds  
12 on mutual information” and our MI estimates are not lower bounds (nor upper bounds). Finally, we thank you for the  
13 relevant additional references.

14 **Reviewer 2:**

15 | *'hard to assess if the claimed severity of the proposed "density-chasm" is valid. ... [the chasm] does not seem to  
16 hinder state-of-the-art results in many applications.'*

17 | *'No comparisons to models other than "single ratio" estimations are given. It raises the question, why this problem  
18 was never addressed before, or if it was addressed before why there is no comparison.'*

19 As noted by reviewers 1 & 3, the “density-chasm problem” is a frequently occurring issue for practitioners using  
20 density-ratio estimation. Despite this fact, we are not aware of any prior work has clearly labelled the issue and provided  
21 a general-purpose solution. This gap in the literature motivated our paper in the first place.

22 In specific applications, it is often possible to design ratio-estimation tasks that avoid the density-chasm problem. As  
23 we show in our MNIST experiments, single-ratio estimation (i.e. NCE) can work very well if the noise distribution is  
24 sufficiently close to the data distribution, as is the case for the RQ-NSF model. A similar strategy (of learning a powerful  
25 noise distribution) has been used many times in the literature, and could be viewed as one of the core motivations  
26 behind GANs. However, learning *both* a ratio-estimator and complex noise distribution, in order to reduce the chasm,  
27 can be challenging/impractical, and hence it is useful for practitioners to have another method at their disposal which  
28 can work with simple noise distributions (as illustrated on MNIST in Figure 5).

29 The strategy of learning a good noise distribution only applies in the context of energy-based modelling. It does not  
30 apply to the more general problem setting of estimating a density-ratio  $p/q$ , where  $p$  and  $q$  are *fixed* in advance. In this  
31 general setting (which includes our MI & representation learning experiments), very few viable approaches exist. We  
32 think that the single density-ratio baselines we compare to—which include results from a 2019 Neurips paper—are  
33 representative of the state-of-the-art.

34 | *'With little analytical guarantees and only few experiments (most of them involving variants of Gaussian noise) it  
35 is hard to assess how generalizable the results are.'*

36 The correctness of TRE is straightforwardly derived from that of single-density ratio estimation, for which there are  
37 extensive analytical guarantees (which we cite in Sec 2). Experiments in 4.1 & 4.2 are intentionally simple illustrations  
38 of the method using synthetic Gaussian data. Experiments in 4.3 & 4.4 use significantly more complex, non-Gaussian  
39 data. Whilst the absolute number of datasets used is modest, we believe the substantial performance increase of TRE  
40 over single ratio methods in a diverse set of applications (MI estimation, representation learning & energy-based  
41 modelling) is strong evidence of its generalizability.

42 | *'The introduced [waymark] mechanisms all appear very ad hoc. How did you decide between the one or another?'*

43 As discussed in the conclusion, we agree further research on the waymark mechanisms would be valuable despite  
44 the good empirical performance of those presented, which were motivated by simplicity. As you noted, across all  
45 applications of TRE we used ‘dim-wise mixing’ for discrete data and the ‘linear combinations’ for continuous data.

46 | *'What did you use as free parameter  $\theta$  (line 88 and Fig.1)? ... [in Fig.1] The axes scaling/cut makes  $p(x)$  invisible...  
47 Can you elaborate on [line 133]? In Fig. 1(a)  $p(x)$  has  $\sigma=1e-6$  and  $q(x)$  has  $\sigma=1$  ... Could you provide a  
48 derivation of Eq. (5) from equal variances (e.g. in the appendix)?'*

49  $\theta$  is given in Sec 3, Eq 7 of the Appendix; we will clarify this.  $p(x)$  is missing in Fig.1 because it perfectly overlaps with  
50 the blue curve; we will amend the caption. Line 133 is slightly wrong: all of our experiments *except* 4.1 (i.e. Fig 1)  
51 preserve variance; we will amend it and add a one-line derivation of the variance property. Thank you for the comments.

52 **Reviewer 4**

53 | *'The main weakness is that the authors never reveal which energy-based model they use in the experiments section.'*

54 We state in Sec 3.4 that the model has the form  $r(\mathbf{x}; \boldsymbol{\theta})q(\mathbf{x})$ , where  $r(\mathbf{x}; \boldsymbol{\theta})$  is the product of all the bridges (Eq. 4) and  
55  $q(\mathbf{x})$  is a noise distribution. In Sec 4.4, we state how each bridge is parameterised, and the different choices of noise  
56 distribution. Sec 1 & 5 of the appendix give all the architectural/training details necessary to reproduce our results.

57 | *'The relationship to AIS could be made a bit more prominent'*

58 We agree, and will include more discussion in the final paper (please see also our response to Revs 1+3).