

1 We thank reviewers for their valuable comments. We will incorporate all feedback in our final manuscript.

2 **Sample Quality/Qualitative/Quantitative Comparisons (R1, R2, R4)** R1/R2/R4 comment on the overall genera-  
3 tion quality of EBMs. We note that the main focus of our paper is to introduce a set of compositional logical operators  
4 over EBMs, and empirical applications, with qualitative fidelity of generations an orthogonal direction. By training an  
5 EBM with a larger number of parameters and computational resources, we see in Figure 1 that EBMs achieve high  
6 fidelity composition results comparable to those of a GAN model (SNGAN is a 128x128 model specifically trained  
7 with Young/Female/Smiling/Wavy Hair attributes). The SNGAN model has the same number of parameters and is  
8 trained with the same number of training iterations using the Mimicry \* GAN library. We will release the code for  
9 training these models and pretrained weights.

10 In terms of image fidelity, on the Young AND Female  
11 AND Smiling AND Wavy Hair split, our composed EBM  
12 obtains an FID of 45.3 while SNGAN obtains FID scores  
13 of 74.2 (all FIDs are large due to a small dataset). R2 fur-  
14 ther asks diversity evaluation. We compare standard deviation  
15 across pixels of generated images and find SNGAN  
16 obtains 55.4 while EBMs obtain 64.5, providing evidence  
17 EBMs generate more diverse samples.

18 **Training Time Comparisons with Other Models (R2)**  
19 A general comparison between training EBMs and other  
20 generative models can be found in the appendix A.5 of [1].  
21 Our EBM models are trained with the same methodology,  
22 and exhibit similar trends. EBMs are slower to train  
23 than GAN models, but faster to train than autoregressive  
24 and flow models. In the particular setting of qualitative  
25 generation in Figure 1, EBM models roughly take 20  
26 times longer to train than the corresponding SNGAN  
27 model (due to generating negative samples).

28 **Relations to Previous Work (R2/R4)** SCAN learns a  
29 fixed latent space for concepts and composition of con-  
30 cepts via logical rules is achieved by manipulating this  
31 latent space. Extending the space of concepts requires retraining the network. With our work, we investigate an  
32 alternative approach where new concepts can be added on demand via new energy functions without invalidating  
33 previous energy functions. This unique characteristic allows for unique benefits – such as the ability to learn visual  
34 concepts in a continual manner. Different from past work in EBMs, our approach is the first to propose additional logical  
35 operators of disjunction and negation, and show that these logical operators can be composed and nested together.

36 **Compositions of Non-Overlapping Concepts (R2)** We train separate EBMs on the frogs  
37 and trucks CIFAR-10 images. We combine models in Figure 2, and find somewhat reasonable  
38 generations that share properties in both classes, although we do not expect good results in  
39 this regime generally.

40 **Comparison to Attention Masks (R2)** While image masking approaches to composi-  
41 tionality enable a part-like decomposition of a scene, generation of each part is largely  
42 independent. This can miss interaction effects between parts (such as shadow casting). Our  
43 EBM composition generates all pixels of the image jointly, offering potential of capturing  
44 such interaction effects.

45 **Equality of Partition Functions (R3)** We had difficulty in checking for equality of parti-  
46 tion functions without using spectral or L2 regularization as they are necessary for stable  
47 training of our method. We will clarify this in the paper.

48 **Multiple Energy Functions (R3)** Ensuring that individual energy functions all have good generative performance  
49 can be difficult. We find that using our proposed operators to compose generation can lessen the need for any individual  
50 model to have good generative performance (see for example Figure 1).

51 **Additional Feedback/Clarification/Typos (R2, R4)** We will add remarks in figure 1, and merge 3.3 with experiments  
52 section and the generalization. We will add a longer introduction about the EBMs and early equations. We will further  
53 fix the typos from R2 and R4, including the Logsumexp expression.

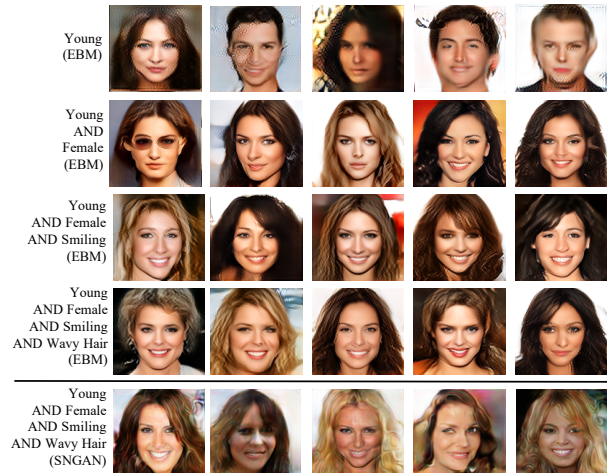


Figure 1: High resolution samples of attribute compositionality with EBMs (same setup as Figure 3 in the main paper). The last row shows SNGAN samples trained on specific attribute combination.

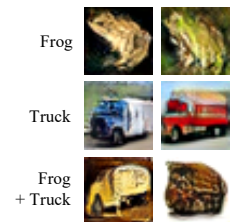


Figure 2: Hybrid combinations of frog and truck EBMs.

\* <https://github.com/kwotsin/mimicry>