

---

# Compositional Visual Generation with Composable Diffusion Models

---

**Nan Liu** \*  
UIUC  
nanliu4@illinois.edu

**Shuang Li** \*†  
MIT CSAIL  
lishuang@mit.edu

**Yilun Du** \*†  
MIT CSAIL  
yilundu@mit.edu

**Antonio Torralba**  
MIT CSAIL  
torralba@mit.edu

**Joshua B. Tenenbaum**  
MIT CSAIL, BCS, CBMM  
jbt@mit.edu

## Abstract

Large text-guided diffusion models, such as DALLÉ-2, are able to generate stunning photorealistic images given natural language descriptions. While such models are highly flexible, they struggle to understand the composition of certain concepts, such as confusing the attributes of different objects or relations between objects. In this paper, we propose an alternative structured approach for compositional generation using diffusion models. An image is generated by composing a set of diffusion models, with each of them modeling a certain component of the image. To do this, we interpret diffusion models as energy-based models in which the data distributions defined by the energy functions may be explicitly combined. The proposed method can generate scenes at test time that are substantially more complex than those seen in training, composing sentence descriptions, object relations, human facial attributes, and even generalizing to new combinations that are rarely seen in the real world. We further illustrate how our approach may be used to compose pre-trained text-guided diffusion models and generate photorealistic images containing all the aspects described in the input descriptions, including the binding of certain object attributes that have been shown difficult for DALLÉ-2. These results point to the effectiveness of the proposed method in promoting structured generalization for visual generation. Project page: <https://energy-based-model.github.io/Compositional-Visual-Generation-with-Composable-Diffusion-Models/>

## 1 Introduction

Our understanding of the world is highly compositional in nature. We are able to rapidly understand new objects from their components, or compose words into complex sentences to describe the world states we encountered [20]. We are able to make ‘infinite use of finite means’ [4], *i.e.*, repeatedly reuse and recombine concepts we have acquired in a potentially infinite manner. We are interested in constructing machine learning systems to have such compositional capabilities, particularly in the context of generative modeling.

Existing text-conditioned diffusion models, such as DALLÉ-2 [32], have recently made remarkable strides towards compositional generation, and are capable in generating photorealistic images given

---

\*indicates equal contribution

†Correspondence to: Shuang Li <lishuang@mit.edu>, Yilun Du <yilundu@mit.edu>