Unsupervised Learning of Compositional Energy Concepts

Yilun Du MIT CSAIL yilundu@mit.edu Shuang Li MIT CSAIL lishuang@mit.edu Yash Sharma
University of Tübingen
yash.sharma@uni-tuebingen.de

Joshua B. Tenenbaum MIT CSAIL, BCS, CBMM jbt@mit.edu Igor Mordatch Google Brain imordatch@google.com

Abstract

Humans are able to rapidly understand scenes by utilizing concepts extracted from prior experience. Such concepts are diverse, and include global scene descriptors, such as the weather or lighting, as well as local scene descriptors, such as the color or size of a particular object. So far, unsupervised discovery of concepts has focused on either modeling the global scene-level or the local object-level factors of variation, but not both. In this work, we propose COMET, which discovers and represents concepts as separate energy functions, enabling us to represent both global concepts as well as objects under a unified framework. COMET discovers energy functions through recomposing the input image, which we find captures independent factors without additional supervision. Sample generation in COMET is formulated as an optimization process on underlying energy functions, enabling us to generate images with permuted and composed concepts. Finally, discovered visual concepts in COMET generalize well, enabling us to compose concepts between separate modalities of images as well as with other concepts discovered by a separate instance of COMET trained on a different dataset*.

1 Introduction

Human intelligence is characterized by its ability to learn new concepts, such as the manipulation of a new tool from only a few demonstrations [1]. Essential to this capability is the composition and re-utilization of previously learned concepts to accomplish the task at hand [38]. This is especially apparent in natural language, which is often described as a tool for making 'infinite use of finite means' [9]. Previously acquired words can be infinitely nested using a set of grammatical rules to communicate an arbitrary thought, opinion, or state one is in. In this work, we are interested in constructing a system that can discover, in an unsupervised manner, a broad set of these compositional components, as well as subsequently combine them across distinct modalities and datasets.

For obtaining such decompositions, two separate lines of work exist. The first focuses on obtaining global, holistic, compositional factors by situating data points, such as human faces, in an underlying (fixed) factored vector space [25, 51]. Individual factors, such as emotion or hair color, are represented as independent dimensions of the vector space, with recombination between factors corresponding to the recombination of the underlying dimensions. Due to the fixed dimensionality of the vector space, multiple instances of a single factor, such as lighting, may not be combined, nor can individual

^{*}Code and data available at https://energy-based-model.github.io/comet/

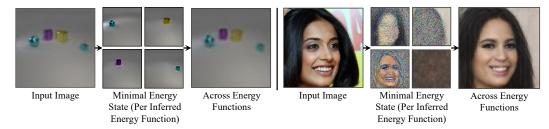


Figure 1: COMET decomposes images into a set of energy functions. The minimal energy state across all energy functions reconstructs the input image. The minimal energy states of individual energy functions capture particular aspects of an image, in the form of local factors of variations such as individual objects, or global factors of variation such as hair color, background lighting, facial expression, or skin color. **Note that reconstructed face images in the remainder of the paper are synthetic and not existing human faces.**

factored vector spaces from separate datasets be combined, i.e. the facial expression in an image from one dataset, and the background lighting in an image from another.

To address this weakness, a separate line of work decomposes a scene into a set of underlying 'object' factors. Each object factor represents a separate set of pixels in an image, as defined by a disjoint segmentation mask [23, 50]. Such a representation allows for the composition of individual factors by compositing the segmentation masks. However, by explicitly constraining the decomposition to be in terms of disjoint segmentation masks, relationships between individual factors of variation become more difficult to represent and capture global factors describing a scene.

In this work, we propose instead to decompose a scene into a set of factors represented as *energy functions*. An individual energy function represents a factor by assigning low energy to scenes with said factor and high energy to scenes where said factor is absent. A scene is then generated by optimizing the sum of energies for all factors. Multiple factors can be composed together, by summing the energies for each individual factor. Simultaneously, these individual factors are defined across the entire scene, allowing energy functions to represent global factors (lighting, camera viewpoint) and local factors (object existence).

Our work is inspired by recent work [13] which shows that energy based models may be utilized to represent flexible compositions of both global and local factors. However, while [13] required supervision, as labels were used to represent concepts, we aim to decompose and discover concepts in an unsupervised manner. In our approach, we discover energy functions from separate data points by enforcing compositions of these energy functions to recompose data.

Our work is further inspired by the ability of humans to effortlessly and reliably combine concepts gathered from disparate experiences. As a separate benefit of our approach, we show that we may take components inferred by one instance of our model trained on one dataset, and compose it with other components inferred by other instances of our model trained on separate datasets.

We provide analysis showing why our approach, COMET[†], is favorable compared to existing unsupervised approaches for decomposing scenes. We contribute the following: First, we show COMET provides a unified framework enabling us to decompose images into both global factors of variation as well as local factors of variation. Second, we show that COMET enables us to scale to more realistic datasets than previous work. Finally, we show that components obtained by COMET generalize well, and are amenable to compositions across different modes of data, and with components discovered by other instances of COMET.

2 Related Work

Our work is related to research in the areas of global factor disentanglement [5, 8, 24, 36, 40, 46] and independent component analysis (ICA) [3, 10, 26–28, 32, 33, 45, 54]. Work in said areas focus on discovering an underlying global latent space which describes the input space. In contrast, we decompose data into a set of compositional vector spaces. This enables our approach to compose multiple instances of one factor together, as well as compose factors across distinct datasets.

[†]short for COMposable Energy neTwork

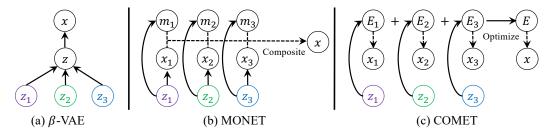


Figure 2: Illustration of distinct approaches to composing components z_1, z_2, z_3 into an image x. (a) β -VAE utilizes a global decoder to map components to images. (b) MONET composits disjoint segmentation masks representing each image. (c) COMET defines an energy function per component, and optimizes the sum over the set of energy functions for each component.

Our work is also related to a large body of existing work on unsupervised object discovery [6, 11, 15–18, 22, 23, 37, 41, 48, 49, 52]. Such methods seek to decompose the scene into its underlying compositional objects by independently segmenting out each individual object in the scene. In contrast, our approach represents individual objects without the need of an explicit segmentation mask, enabling our approach to represent global relations between objects.

Our work draws on recent work in energy based models (EBMs) [12, 14, 19, 21, 34, 43, 47, 53]. Our underlying energy optimization procedure to generate samples is reminiscent of Langevin sampling, which is used to sample from EBMs [12, 43, 53]. Most similar to our work is that of [13], which proposes composing EBMs for compositional visual generation. Different from [13], we study how we may discover factors in an unsupervised manner from data.

3 Composable Energy Networks

Let $\mathcal{D}=\{X\}$ be a set of images $x\in\mathbb{R}^D$. Our goal is to obtain a set Z of components $\{z_0,z_1,\ldots,z_k\}$ for each image x representing the underlying factors of variation in the image, where each component $z\in\mathbb{R}^M$. We first discuss how we represent Z as a set of energy functions in Section 3.1. We then discuss how COMET learns to decompose data in an unsupervised manner into a set of components in Section 3.2.

3.1 Composing Factors of Variation as Energy Functions

In prior work for decomposing images into a set of components Z, approaches such as β -VAE [24] utilize a parameterized, non-adaptive, feedforward decoder to map the set of components to an image. Due to this, the method is unable to compose multiple instances of the same component or a larger number of components then seen during training. On the other hand, approaches such as MONET [6] utilize a decoder with shared weights to process each individual component. Unfortunately, the components are decoded independently, preventing relationships between individual components from being captured.

To flexibly and generically compose a set of components Z, we require an approach that (1) utilizes a decoder that is shared across each component, such that variable-sized sets of components may be encoded, while also (2) decodes components jointly such that the relationships between individual components are covered. To construct an approach with the aforementioned desiderata, we propose to represent each component z as an energy function. We illustrate our approach and its differences from prior work in Figure 2.

Representing Components as Energy Functions. Given a single component z, we encode the factor using an energy function $E_{\theta}(x,z): \mathbb{R}^D \times \mathbb{R}^M \to \mathbb{R}$, which is learned to assign low energy to all images x which contain the component and high energy to all other images. To obtain an image x' with a component z, we solve for the expression $x' = \arg\min_{x} E_{\theta}(x; z_i)$. Note that in contrast to prior work composing energy functions [13], our energy functions have no probabilistic interpretation.

Composing Energy Functions. Given multiple components, we represent a set of components $\{z_i\}$ by performing a summation $\sum_i E_{\theta}(x, z_i)$ over each component's energy function. To obtain an image x' for the set of components Z, we solve for the expression $x' = \arg\min \sum_i E_{\theta}(x; z_i)$.

3

While both COMET and MONET rely on summation as a tool for composing components, MONET sums components in the image domain, while COMET sums the cost functions representing each component together. By summing cost functions, each component may combine with other components in a non-linear manner, enabling us to model more complex relationships between individual components.

As a result, our generation x' is the byproduct of jointly minimizing each individual energy function $E_{\theta}(x, z_i)$, and thus our generation contains each component. In addition, each individual component z_i is parameterized by the same energy function E_{θ} , enabling us to model a different number of components.

3.2 Unsupervised Decomposition of Composable Energy Functions

We next discuss how COMET discovers a set of composable energy functions from an input image x_i . In Section 3.1, we discuss that a set of components $Z = \{z_i\}$ may be encoded as a set of composable energy functions using the expression $\arg\min_{\boldsymbol{x}}\sum_i E_{\theta}(\boldsymbol{x}; z_i)$. To discover the set of components Z in an unsupervised manner, we train by recomposing the input image x_i

$$\mathcal{L}_{\text{MSE}}(\theta) = \|\arg\min_{\boldsymbol{x}} \left(\sum_{k} E_{\theta}(\boldsymbol{x}; \text{Enc}_{\theta}(\boldsymbol{x}_i)[k])\right) - \boldsymbol{x}_i\|^2.$$
 (1)

We utilize a learned neural network encoder $\operatorname{Enc}_{\theta}(\boldsymbol{x}_i)$ to infer a set of components \boldsymbol{z}_k . In practice, evaluating the expression $\operatorname{arg\,min}_{\boldsymbol{x}} \sum_k E_{\theta}(\boldsymbol{x}; \operatorname{Enc}_{\theta}(\boldsymbol{x}_i)[k])$, is computationally intractable. We thus approximate the $\operatorname{arg\,min}$ operation with respect to \boldsymbol{x} via N steps of gradient optimization, with an approximate optimum \boldsymbol{x}_i^N obtained as

$$\boldsymbol{x}_{i}^{N} = \boldsymbol{x}_{i}^{N-1} - \lambda \nabla_{\boldsymbol{x}} \sum_{k} E_{\theta}(\boldsymbol{x}_{i}^{N-1}; \operatorname{Enc}_{\theta}(\boldsymbol{x}_{i})[k]). \tag{2}$$

In the above expression, we initialize optimization of \boldsymbol{x}_i^0 with uniform noise with λ representing the step size for each gradient step. Sample \boldsymbol{x}_i^n denotes the result after n steps of gradient descent. We train our energy function E_{θ} using the modified objective $\mathcal{L}_{\text{MSE}} = \sum_{n=1}^N \|\boldsymbol{x}_i^n - \boldsymbol{x}_i\|^2$, where we minimizing MSE w.r.t \boldsymbol{x}_i^n . To train parameters of E_{θ} we use automatic differentiation to compute gradients with respect to each optimization step depicted in Equation 2, where for training stability, we truncate gradient backpropogation to the previous time step.

We provide pseudocode for training our model in Algorithm 1. While the approach is simple, we find that it performs favorably in both global disentanglement as well as object-level disentanglement, and requires no additional objectives Algorithm 1 Training algorithm for COMET.

Input: data dist. $p_D(x)$, step size λ , number of gradient steps N, encoder Enc_{θ} , energy function E_{θ} , energy components K

while not converged do

```
\begin{aligned} & \boldsymbol{x}_i \sim p_D \\ & \triangleright \textit{Encode components } \boldsymbol{z}_i^k \textit{ from } \boldsymbol{x}_i \\ & \boldsymbol{z}_i^1, \dots, \boldsymbol{z}_i^K \leftarrow \text{Enc}_{\theta}(\boldsymbol{x}_i) \\ & \triangleright \textit{Optimize sample } \boldsymbol{x}_i^0 \textit{ via gradient descent:} \\ & \boldsymbol{x}_i^0 \sim \mathcal{U}(0,1) \\ & \textbf{for } \text{gradient step } n = 1 \text{ to } N \textbf{ do} \\ & \boldsymbol{x}_i^n \leftarrow \boldsymbol{x}_i^{n-1} - \lambda \nabla_{\boldsymbol{x}} \sum_{k=1}^K E_{\theta}(\boldsymbol{x}_i^{n-1}; \boldsymbol{z}_i^k) \\ & \textbf{end for} \\ & \triangleright \textit{Optimize objective } \mathcal{L}_{\textit{MSE}} \textit{ wrt } \theta \text{:} \\ & \Delta \theta \leftarrow \nabla_{\theta} \sum_{n=1}^N \|\boldsymbol{x}_i^n - \boldsymbol{x}_i\|^2 \\ & \text{Update } \theta \textit{ based on } \Delta \theta \textit{ using optimizer} \\ & \textbf{end while} \end{aligned}
```

or priors to shape underlying latents. We utilize the same energy architecture E_{θ} throughout all experiments, and present details in the appendix.

Controlling Local and Global Decompositions. In many scenes, both global and local factor decompositions are valid. To control the decomposition obtained by COMET, we bias the system towards inferring local factor decompositions by utilizing i) low latent dimensionality and ii) positional embeddings [39], both of which have been used in previous object discovery works [6]. This bias serves to encourage and enable models to focus on local patches of an image. We provide analysis of the effect of each inductive bias on the decomposition in Section 5.2.

4 Complexity Analysis of Energy Function Compositions

In this section, we provide complexity-theoretic motivation for composing functions together using energy functions. Let Z denote a set of components, with individual components $z_1, \ldots, z_K \in \mathbb{R}^M$.

Let \mathcal{F} be the space of functions $f: Z \to \mathbb{R}^D$ for mapping sets of components to an output vector x. We consider two subspaces for \mathcal{F} .

Composition of Energy Functions. Let $\mathcal{F}_{\text{energy}} \subseteq \mathcal{F}$ be the subspace of functions of the form $\arg\min_{\boldsymbol{x}} \sum_{1 \leq k \leq K} E(\boldsymbol{x}, \boldsymbol{z}_k)$, where each $E(\boldsymbol{x}, \boldsymbol{z}_m)$ is a function from $\mathbb{R}^D \times \mathbb{R}^M \to \mathbb{R}$. This subspace corresponds to the set of compositions realized by COMET.

Composition of Segmentation Masks. We next consider the subspace of functions $\mathcal{F}_{\text{mask}} \subseteq \mathcal{F}$ consisting of functions f of the form $\sum_{1 \leq k \leq K} m_k(\boldsymbol{z}_k) f_k(\boldsymbol{z}_k)$ where $f(\boldsymbol{z}) : \mathbb{R}^M \to \mathbb{R}^D$ and $m_k(\boldsymbol{z}_k) : \mathbb{R}^M \to \{0,1\}^D$. Here, $m_k(\boldsymbol{z}_k)$ represents a segmentation mask in \mathbb{R}^D . This composition is used in object decomposition methods, such as [6,23].

We first show that compositions of energy functions are more expressive than compositions of segmentation masks.

Remark 1. The subspace \mathcal{F}_{energy} is a strict superset of the subspace \mathcal{F}_{mask} .

Proof. Any function of the form $\sum_{z \in Z} m_k(z) f(z)$ is equivalent to $\arg\min_x \sum_{z \in Z} E'(x,z)$, where $E'(x,z) = m_k(z)(x-f(z))^2$, so $\mathcal{F}_{\text{mask}} \subseteq \mathcal{F}_{\text{energy}}$. In the other direction, note that according to the definition of $\mathcal{F}_{\text{mask}}$, the presence of a particular component z_i assigns a fixed value to the output for the nonzero entries of $m_k(z_k)f(z_k)$, irrespective of the value of the other components. In contrast, the optimal value of x in $\mathcal{F}_{\text{energy}}$ depends on the value of all components. A constructive example of this is the set of energy functions over one-dimensional x where $E(x,z_1)=(x-2)^2$, $E(x,z_2)=(x-3)^2$, $E(x,z_3)=(x-4)^2$. Thus, we have that $\arg\min_x E(x,z_1)+E(x,z_2)\neq\arg\min_x E(x,z_1)+E(x,z_3)$. Therefore, there are functions in $\mathcal{F}_{\text{energy}}$ that are not in $\mathcal{F}_{\text{mask}}$.

Next, we show that even in the setting in which we represent a single component z, learning a function to approximate E(x,z) is more computationally efficient than learning a function f(z) that approximates $\arg\min_{\boldsymbol{x}} E(\boldsymbol{x},z)$. Intuitively, an energy function $E(\boldsymbol{x},z)$ can be seen as a verifier of a set of constraints, with the value of $E(\boldsymbol{x},z)$ being low when all constraints are satisfied. Approximating the function $\arg\min_{\boldsymbol{x}} E(\boldsymbol{x},z)$ corresponds to generating a solution given a set of constraints, which is well-known in complexity theory to be much harder than verifying the constraints. Thus, to enable formal analysis of $E(\boldsymbol{x},z)$, we reduce the 3-SAT [29] problem to an energy function $E(\boldsymbol{x},z)$.

Given a 3-SAT formula ϕ with D variables and K clauses, we encode ϕ using an energy function $E(\boldsymbol{x},\boldsymbol{z}) := \sum_{1 \leq k \leq K} e_k(\boldsymbol{x},\boldsymbol{z}[k])$, where, \boldsymbol{x} encodes an assignment to all the variables, $\boldsymbol{z}[k]$ represents the dimensions of \boldsymbol{z} encoding the k^{th} clause, and each e_k is a function which has energy 0 if the assignment \boldsymbol{x} satisfies the k^{th} clause, and has energy 1 otherwise. To encode a clause using $\boldsymbol{z}[k]$, we utilize an ordinal representation (e.g. $\boldsymbol{z} = [1, 2, 3]$ to represent the clause $(x_1 \wedge x_2 \wedge x_3)$), and round non-integer coordinates of \boldsymbol{x} and \boldsymbol{z} to the nearest integer. We assume the Exponential Time Hypothesis (ETH)[29], which states that checking the satisfiability of a 3-SAT formula takes time exponential in the sum of the number of variables and the number of clauses.

Remark 2. There exists an energy function $E(x, z_1)$ which can be evaluated at any input in time polynomial in the number of dimensions of x but for which the computational complexity of evaluating $f(z) := \arg\min_{x} E(x, z_1)$ is (worst-case) exponential in the number of dimensions of x.

Proof. If we utilize the 3-SAT energy function E(x, z) defined above. ETH implies that solving the 3-SAT problem, corresponding to computing $f(z) = \arg\min_{x} E(x, z_1)$, is exponential in dimension of x. In contrast, evaluating each entry of our defined $E(x, z_1)$ is polynomial in dimension of x. \square

Our remark shows that it is computationally advantageous to learn an energy function E(x, z) as opposed to a decoder f(z). To realize the exponential number of computations needed to compute f(z), significantly more capacity is necessary to represent f(z) in comparision to E(x, z). We further show that as we compose multiple 3-SAT energy functions together, learning a decoder $f(z_1, \ldots, z_k) := \arg\min_x \sum_k E(x, z_k)$ that represents the optimization process scales exponentially with the number of energy components.

Remark 3. There exists a composition of energy functions $\sum_k E(\boldsymbol{x}, \boldsymbol{z}_k)$ which can be evaluated in time polynomial in the number of components k, but for which the computational complexity for evaluating $f(\boldsymbol{z}_1, \dots, \boldsymbol{z}_k) := \arg\min_{\boldsymbol{x}} \sum_k E(\boldsymbol{x}, \boldsymbol{z}_k)$ is (worst-case) exponential in the components of components k.

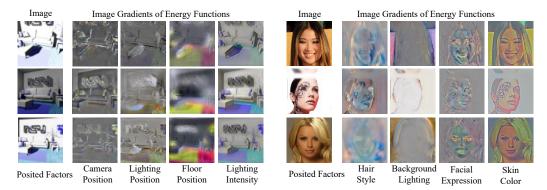


Figure 3: **Global Decomposition.** Illustration of energy functions gradients of each decomposed energy function in COMET on Falcor3D (**left**) and CelebA-HQ (**right**) datasets. Gradients correspond to aspects of images each energy function pays attention to. Discovered energy functions are labeled with the posited factors they capture, as determined by qualitative inspection.

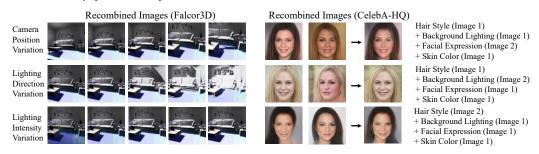


Figure 4: **Global Factor Recombination.** Illustration of recombination of energy functions on Falcor3D and CelebA-HQ datasets. In Falcor3D (**left**), we illustrate variation of a single energy function, which elicits changes in camera position, lighting direction and lighting intensity variation. In CelebA-HQ (**right**), we recombine discovered energy functions across two separate images (energy functions labeled through qualitative visualization in Figure 3).

Proof. Given k separate 3-SAT energy functions $E(\boldsymbol{x}, \boldsymbol{z}_k)$, minimizing the composed energy function, $f(\boldsymbol{z}_1, \dots, \boldsymbol{z}_k) = \arg\min_{\boldsymbol{x}} \sum_k E(\boldsymbol{x}, \boldsymbol{z}_k)$ corresponds to solving all 3-SAT encoded clauses across k energy functions. ETH implies the computational complexity of evaluating $f(\boldsymbol{z}_1, \dots, \boldsymbol{z}_k)$ is exponential in k while the evaluation of k energy functions is polynomial in k.

Here as well, to represent the exponential number of computations, significantly more capacity is necessary to realize the computation $f(z_1, \ldots, z_k)$, which scales with the number of components k, in comparison to $\sum_k E(x, z_k)$. Thus, remark 2 and 3 show that it can be efficient, from a learning perspective, to decode individual datapoints with energy functions.

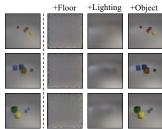
5 Evaluation

We quantitatively and qualitatively show that COMET can recover the global factors of variation in an image in Section 5.1, as well as the local factors in an image in Section 5.2. Furthermore, we show that the components captured by COMET can generalize well, across separate modalities in Section 5.3.

5.1 Global Factor Disentanglement

We assess the ability of COMET to decompose global factors of variation in scenes consisting of lighting and camera illumination from Falcor3D (NVIDIA high-resolution disentanglement dataset) [42], scene factors of variation in CLEVR [30], and face attributes in real images from CelebA-HQ [31]. For all experiments, we utilize the same convolutional encoder to extract sets of latents from each dataset, and utilize a latent dimension of 64 for each separately inferred latent. We provide additional training algorithm and model architecture details in the appendix.

Decomposition and Reconstructions. In Figure 3, on Falcor3D and CelebA-HQ, we illustrate the underlying image gradients of each decomposed energy function. Such gradients correspond to



Input Generation (Per Additive Component)

Figure 5: **CLEVR Decomposition.** Generations on CLEVR when optimizing over an increasing number of energy functions.

Model	Dim (D)	β	Decoder Dist.	BetaVAE	MIG	MCC
COMET* β-VAE	64 64	- 4	– Gaussian	99.41 ± 0.15 83.57 ± 8.05	19.63 ± 2.49 10.90 ± 3.80	76.55 ± 1.35 66.08 ± 2.00
β-VAE β-VAE	32 256	4 4	Gaussian Gaussian	79.77 ± 10.95 80.76 ± 4.55	7.14 ± 5.44 10.94 ± 0.58	57.48 ± 6.04 66.14 ± 1.81
β-VAE β-VAE	64 64	16 1	Gaussian Gaussian	74.71 ± 1.57 81.61 ± 6.75	9.33 ± 3.72 6.51 ± 3.38	57.28 ± 2.37 58.73 ± 6.31
β-VAE	64	4	Bernoulli	84.23 ± 3.51	8.96 ± 3.53	61.57 ± 4.09
InfoGAN MONet*	64 64	_		79.65 ± 1.69 93.13 ± 1.02	2.48 ± 1.11 13.94 ± 2.09	52.67 ± 1.91 65.72 ± 0.89

Table 1: **Disentanglement Evaluation.** Mean and standard deviation (s.d.) metric scores across 3 random seeds on the Falcor3D dataset. COMET enables better disentanglement according to 3 common disentanglement metrics across different runs and seeds for training β -VAE, InfoGAN and MONet baselines. Note that * denotes that PCA was used as a postprocessing step.

aspects of the input image each energy function pays attention. Through qualitative examination, we posit that the individual inferred energy functions for Falcor3D correspond to camera position, lighting direction, floor position and lighting intensity (shown from left to right in Figure 3). Such a correspondence can be seen for example by the fact that the camera position energy function exhibits sharp gradients with respect to edges in an image, while the lighting direction energy function exhibits sharp gradients with respect to the underlying shadows in an image. In a similar manner, we hypothesize that the individual inferred energy functions for CelebA-HQ correspond to hair color, background lighting, facial expression and skin color.

In Figure 5, we show energy function decompositions of CLEVR, where we illustrate generations when composing an increasing number of inferred energy functions. By adding individual energy functions, we find that the generation transitions from consisting only of objects, to consisting of objects & floor to consisting of objects, floor, & lighting.

Recombination. To further verify that each individual energy function captures the expected decomposition of factors described earlier, we recombine individual energy functions representing each component in the Falcor3D and CelebA-HQ datasets in Figure 4. In the left side of Figure 4, we vary an energy function representing a single factor, while keeping the remaining factors fixed. By varying energy functions in such a manner, we are able to capture camera position, lighting direction and lighting intensity variation (with camera position variation captured by varying energy functions for both floor position and camera position). In the right side of Figure 4, we can recombine discovered energy functions representing facial expression, hair color, and background lighting of one image with that of another. We find that by recombining individual energy functions representing each individual factor, we are able to reliably swap factors across images.

Quantitative Comparison. Finally, we evaluate the learned representations on disentanglement. In Falcor3D [42], each image corresponds to a combination of 7 factors of variation; lighting intensity, lighting x, y & z direction, and camera x, y & z position. We consider three commonly used metrics for evaluation, the BetaVAE metric [24], the Mutual Information Gap (MIG) [8], and the Mean Correlation Coefficient (MCC) [26]. See Appendix C of [36] for extended descriptions of metrics.

Standardized metrics in disentanglement assume flattened model latents, i.e. relate the 7 factors of variation to the D-dimensional encoding of the corresponding image. However, as discussed, our method extracts sets of latents from images. We thus extract D principal components from the set of latents, which we then compare with β -VAE. In Table 1, we find that our approach performs better than that of β -VAE across hyperparameter settings, as well as additional baselines of MONet and InfoGAN.

5.2 Object Level Decomposition

Next, we assess the ability of COMET to decompose object-level factors of variation in an image. We evaluate the ability of COMET to isolate individual tetris blocks in the Tetris dataset from [23], and individual object segmentations in the CLEVR dataset [30]. To bias energy functions to local object-level variations, we utilize small latent dimension (16), and add positional embeddings to images. To extract repeated object-level structure in an image, we utilize a recurrent network as our

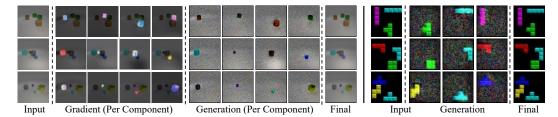


Figure 6: **Object Decomposition.** COMET can decompose underlying object-level factors of variation in CLEVR (left) and Tetris (right). Individual energy functions exhibit sharp gradients with respect to each object in CLEVR.

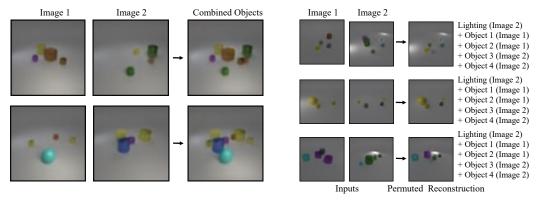


Figure 7: **Compositional Generalization (Left).** We may combine energy functions from COMET on two separate images together to generate a new image with 8 objects. **Global/Local Factor Recombination (Right).** COMET is able to discover energy functions corresponding to global and local factors of variation simultaneously from a given image. We can recombine individual energy functions corresponding to global factors of variation (lighting) and local factors of variation (cube position).

encoder, similar to prior unsupervised object discovery work [6]. We provide additional training algorithm and model architecture details for each experimental setup in the appendix.

Decomposition and Reconstructions. We illustrate object-level decomposition of individual energy functions on CLEVR and Tetris in Figure 6. In both settings, we find that individual energy functions correspond to the underlying objects in the scene. On CLEVR, while in Figure 5 we obtain global factor decompositions, by biasing our energy function to object level decompositions we obtain individual cube decompositions in Figure 6.

Combinatorical Recombination. We next consider recombining inferred object components. In Figure 7 (left), we combine energy functions for individual CLEVR objects across two separate scenes. We find that by combining 8 energy functions corresponding to 4 objects from separate images, we are able to successfully generate an image containing all 8 objects, and with some consistency with respect to occlusion. In contrast, objects composed together through MONET do not respect occlusion, as they are represented as disjoint segmentation masks.

Quantitative Comparisons. For quantitative comparison with MONET, we create approximate segmentation masks per energy function in CLEVR by thresholding the gradients (illustrated in Figure 6) of the energy function. We find that our approach obtains an ARI [23] of 0.916 and a mean segmentation covering of 0.713. In contrast, we find MONET obtains an ARI of 0.873 and a mean segmentation covering of 0.701.

Decompositions of Global and Object Level Factors. COMET is distinct from previous approaches for unsupervised decomposition in that it is able to decompose both global and local factors of variation. We find that we are further able to *control* the capture of either a particular global or local factor. In particular, we find that by conditioning a particular energy function with a positional embedding added to the image [39] and a small latent dimension we may effectively bias a energy function to capture an object factor of the scene. In contrast, a larger latent size biases an energy function to capture a global factor of the scene. In Table 2, we investigate the extent of these two effects in enabling the capture of local factors in CLEVR as measured by ARI. We find that the

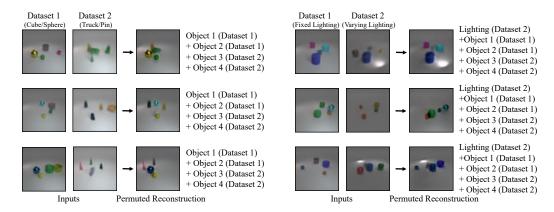


Figure 8: **Cross Dataset Composition.** Energy functions from COMET discovered in one dataset may be combined with energy functions from a separate instance of COMET discovered on a different dataset. In the **left panel**, we recombine energy functions representing cube/sphere objects discovered in CLEVR with energy functions representing truck/pin objects discovered in CLEVR Toy. We are able to generate novel images consisting of objects from both datasets. In the **right panel**, we recombine energy functions representing separate objects in CLEVR with energy functions representing both lighting and individual objects from CLEVR Lighting. We are able to generate novel images consisting of novel combinations of both objects and lighting.

addition of utilizing both small latent dimension and positional embedding enable more effective decomposition of underlying object factors.

By selectively enabling ourselves to specify both the underlying global and local factors of a scene, we may obtain an unsupervised decomposition of a scene with both local and global factors of variation. To test this, we render a novel dataset, *CLEVR Lighting*, by increasing the lighting variation in CLEVR (illustrated in Figure 7 (right)), and train COMET to infer 5 separate energy functions. We encourage the first energy function to capture a global factor of variation by removing the positional embedding from the first inferred energy function. In such a setting, we find COMET successfully infers lighting as the first energy function, and individual objects for the subsequent components, which we illustrate in the appendix. We

Small	Pos	ARI ↑
Latent	Embed	
No	No	0.413
No	No	0.407
Yes	No	0.641
Yes	Yes	0.889
	No No Yes	LatentEmbedNoNoNoNoYesNo

Table 2: Local Factor Inductive Bias. Analysis of different inductive biases on underlying local factor decomposition as measured by ARI score on the CLEVR dataset.

present recombinations of these inferred components in Figure 7 (right). A limitation of our approach is that that lighting and object factors of variations are not completely disentangled. For example, in the top row, the permuted reconstruction of the image has incorrect lighting in the center.

5.3 Compositional Factor Generalization

Finally, we assess the ability of components inferred from COMET to generalize. We study two separate settings of generalization. We first evaluate the ability of components to generalize to multi-modal inputs by training COMET to decompose images drawn from both CelebA-HQ and Danbooru datasets [4], as well as KITTI [20] and Virtual KITTI [7] datasets. We next assess the ability of components from COMET to compose with a separate instance of COMET. We train separate COMET models on CLEVR, CLEVR Lighting and an additional novel dataset, *CLEVR Toy*, which is rendered by replacing sphere/cylinder/cube objects with pin/boot/toy/truck objects.

Cross Dataset Recombination. COMET may compose inferred components from one dataset with components discovered by a separate COMET trained on a separate dataset. We validate this in Figure 8 (left), where we recombine components specifying objects in CLEVR and CLEVR Toy. In Figure 8 (right), we further recombine components specifying objects in CLEVR with components specifying objects & lighting in CLEVR Lighting.

Cross Modality Decomposition and Recombination. We assess COMET decompositions on multi-modal datasets of Danbooru/CelebA-HQ and KITTI/Virtual KITTI in Figure 9. We find consistent decomposition of components between separate modes. Furthermore, we find that these components may be recombined between the separate modalities in Figure 10.

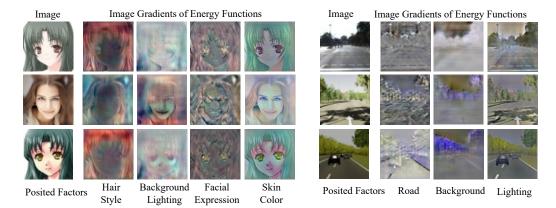


Figure 9: **Multi-modal Decomposition.** Illustration of energy gradients of each unsupervised decomposed energy function in COMET on sets of images drawn from distinct modalities. COMET discovers consistent decompositions between different modalities, where energy functions are labeled with the posited factor they capture (as determined from qualitative inspection). (**Left**) Visualization on images in the Danbooru and CelebA-HQ domains. D Energy functions are consistent with those discovered in Figure 3 (CelebA-HQ). (**Right**) Visualization on images in the KITTI and Virtual KITTI domains.

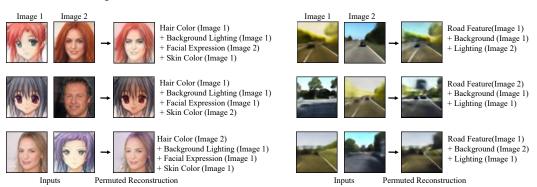


Figure 10: **Multi-modal Recombination.** Illustration of recombination of energy functions on multi-modal datasets of images. On images consisting of Danbooru and CelebA-HQ faces (**left**), by recombining discovered energy functions, we are able to selectively change underlying facial expression, skin color and hair color across images from different domains. On images consisting of KITTI and Virtual KITTI images (**right**), by recombining discovered energy functions, we are able to selectively change the underlying lighting, road and background between images of separate domains.

6 Conclusion

We have demonstrated an approach towards unsupervised learning of energy functions from images. We show how these functions encode both global and local factors of variations, and how they allow for further composition across separate modalities and datasets. A limitation of our current approach is that while it can compose factors across datasets that are substantially similar, compositions across datasets such as Falcor3D and CelebA-HQ are less interpretable. We posit that this is due to a lack of diversity in the underlying training data set. An interesting direction for future work would be to to train COMET on complex, higher diversity real world datasets, and observe subsequent recombinations of energy functions. We note that, as a consequence of using deep nets, our system is susceptible to dataset bias. Care must be put in ensuring a balanced and fair dataset if COMET is deployed in practice, as it should not serve to, even inadvertently, worsen societal prejudice.

Acknowledgements We would like to thank Abhijit Mudigonda for giving helpful comments on the manuscript. Yilun Du is supported by a NSF graduate research fellowship. This work is in part supported by ONR MURI N00014-18-1-2846 and IBM Thomas J. Watson Research Center CW3031624. This work was supported by the German Federal Ministry of Education and Research (BMBF): Tübingen AI Center, FKZ: 01IS18039A. The authors thank the International Max Planck Research School for Intelligent Systems (IMPRS-IS) for supporting Yash Sharma.

References

- [1] Kelsey R. Allen, Kevin A. Smith, and Joshua B. Tenenbaum. Rapid trial-and-error learning with simulation supports flexible tool use and physical reasoning. *Proceedings of the National Academy of Sciences*, 117 (47):29302–29310, 2020. ISSN 0027-8424. doi: 10.1073/pnas.1912341117. URL https://www.pnas.org/content/117/47/29302. 1
- [2] Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio. Neural machine translation by jointly learning to align and translate. In *ICLR*, 2015. 15
- [3] Anthony J Bell and Terrence J Sejnowski. An information-maximization approach to blind separation and blind deconvolution. *Neural Computation*, 7(6):1129–1159, 1995. 2
- [4] Gwern Branwen. Danbooru2019 portraits: A large-scale anime head illustration dataset, 2019. 9
- [5] Christopher P Burgess, Irina Higgins, Arka Pal, Loic Matthey, Nick Watters, Guillaume Desjardins, and Alexander Lerchner. Understanding disentangling in beta-vae. arXiv preprint arXiv:1804.03599, 2018.
- [6] Christopher P Burgess, Loic Matthey, Nicholas Watters, Rishabh Kabra, Irina Higgins, Matt Botvinick, and Alexander Lerchner. Monet: Unsupervised scene decomposition and representation. arXiv:1901.11390, 2019. 3, 4, 5, 8
- [7] Yohann Cabon, Naila Murray, and Martin Humenberger. Virtual kitti 2, 2020. 9
- [8] Tian Qi Chen, Xuechen Li, Roger Grosse, and David Duvenaud. Isolating sources of disentanglement in variational autoencoders. *arXiv:1802.04942*, 2018. 2, 7
- [9] Noam Chomsky. Aspects of the Theory of Syntax. The MIT Press, Cambridge, 1965. URL http://www.amazon.com/Aspects-Theory-Syntax-Noam-Chomsky/dp/0262530074. 1
- [10] Pierre Comon. Independent component analysis, a new concept? *Signal processing*, 36(3):287–314, 1994.
- [11] Eric Crawford and Joelle Pineau. Spatiial invariant unsupervised object detection with convolutional neural networks. In *Thirty-Third AAAI Conference on Artificial Intelligence*, 2019. 3
- [12] Yilun Du and Igor Mordatch. Implicit generation and generalization in energy-based models. *arXiv* preprint arXiv:1903.08689, 2019. 3
- [13] Yilun Du, Shuang Li, and Igor Mordatch. Compositional visual generation with energy based models. In *Advances in Neural Information Processing Systems*, 2020. 2, 3
- [14] Yilun Du, Shuang Li, Joshua B Tenenbaum, and igor Mordatch. Improved contrastive divergence training of energy based models. In *Proceedings of the 38th international conference on Machine learning*. ACM, 2021. 3
- [15] Yilun Du, Kevin A. Smith, Tomer Ullman, Joshua B. Tenenbaum, and Jiajun Wu. Unsupervised discovery of 3d physical objects. In *International Conference on Learning Representations*, 2021. URL https://openreview.net/forum?id=lf7st0bJIA5. 3
- [16] Martin Engelcke, Adam R. Kosiorek, Oiwi Parker Jones, and Ingmar Posner. Genesis: Generative scene inference and sampling with object-centric latent representations. In *International Conference on Learning Representations*, 2020. URL https://openreview.net/forum?id=BkxfaTVFwH.
- [17] SM Eslami, Nicolas Heess, Theophane Weber, Yuval Tassa, Koray Kavukcuoglu, and Geoffrey E Hinton. Attend, infer, repeat: Fast scene understanding with generative models. In *NeurIPS*, 2016.
- [18] SM Ali Eslami, Danilo Jimenez Rezende, Frederic Besse, Fabio Viola, Ari S Morcos, Marta Garnelo, Avraham Ruderman, Andrei A Rusu, Ivo Danihelka, Karol Gregor, et al. Neural scene representation and rendering. *Science*, 360(6394):1204–1210, 2018. 3
- [19] Ruiqi Gao, Erik Nijkamp, Diederik P Kingma, Zhen Xu, Andrew M Dai, and Ying Nian Wu. Flow contrastive estimation of energy-based models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7518–7528, 2020. 3
- [20] Andreas Geiger, Philip Lenz, and Raquel Urtasun. Are we ready for autonomous driving? the kitti vision benchmark suite. In CVPR, 2012. 9
- [21] Will Grathwohl, Kuan-Chieh Wang, Jörn-Henrik Jacobsen, David Duvenaud, Mohammad Norouzi, and Kevin Swersky. Your classifier is secretly an energy based model and you should treat it like one. arXiv preprint arXiv:1912.03263, 2019. 3
- [22] Klaus Greff, Sjoerd van Steenkiste, and Jürgen Schmidhuber. Neural expectation maximization. In NeurIPS, 2017. 3
- [23] Klaus Greff, Raphaël Lopez Kaufmann, Rishab Kabra, Nick Watters, Chris Burgess, Daniel Zoran, Loic Matthey, Matthew Botvinick, and Alexander Lerchner. Multi-object representation learning with iterative variational inference. arXiv preprint arXiv:1903.00450, 2019. 2, 3, 5, 7, 8

- [24] Irina Higgins, Loic Matthey, Arka Pal, Christopher P Burgess, Xavier Glorot, Matthew Botvinick, Shakir Mohamed, and Alexander Lerchner. Beta-vae: Learning basic visual concepts with a constrained variational framework. In *ICLR*, 2017. 2, 3, 7
- [25] Irina Higgins, Nicolas Sonnerat, Loic Matthey, Arka Pal, Christopher P Burgess, Matko Bosnjak, Murray Shanahan, Matthew Botvinick, Demis Hassabis, and Alexander Lerchner. Scan: Learning hierarchical compositional visual concepts. *ICLR*, 2018. 1
- [26] Aapo Hyvärinen and Hiroshi Morioka. Unsupervised feature extraction by time-contrastive learning and nonlinear ica. In Advances in Neural Information Processing Systems, pages 3765–3773, 2016. 2, 7
- [27] Aapo Hyvärinen and Hiroshi Morioka. Nonlinear ica of temporally dependent stationary sources. In Proceedings of Machine Learning Research, 2017.
- [28] Aapo Hyvärinen, Hiroaki Sasaki, and Richard E Turner. Nonlinear ica using auxiliary variables and generalized contrastive learning. arXiv preprint arXiv:1805.08651, 2018.
- [29] R. Impagliazzo and R. Paturi. Complexity of k-sat. Proceedings. Fourteenth Annual IEEE Conference on Computational Complexity (Formerly: Structure in Complexity Theory Conference) (Cat.No.99CB36317), Jun 1999. doi: 10.1109/ccc.1999.766282.
- [30] Justin Johnson, Bharath Hariharan, Laurens van der Maaten, Li Fei-Fei, C Lawrence Zitnick, and Ross Girshick. Clevr: A diagnostic dataset for compositional language and elementary visual reasoning. In CVPR, 2017. 6, 7, 14, 15
- [31] Tero Karras, Timo Aila, Samuli Laine, and Jaakko Lehtinen. Progressive growing of gans for improved quality, stability, and variation. In ICLR, 2017. 6
- [32] Ilyes Khemakhem, Diederik P Kingma, and Aapo Hyvärinen. Variational autoencoders and nonlinear ica: A unifying framework. *International Conference on Artificial Intelligence and Statistics (AISTATS)*, 2020.
 2
- [33] Ilyes Khemakhem, Ricardo Monti, Diederik Kingma, and Aapo Hyvarinen. Ice-beem: Identifiable conditional energy-based deep models based on nonlinear ica. Advances in Neural Information Processing Systems, 33, 2020.
- [34] Taesup Kim and Yoshua Bengio. Deep directed generative models with energy-based probability estimation. *arXiv preprint arXiv:1606.03439*, 2016. 3
- [35] Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In ICLR, 2015. 15
- [36] David A. Klindt, Lukas Schott, Yash Sharma, Ivan Ustyuzhaninov, Wieland Brendel, Matthias Bethge, and Dylan Paiton. Towards nonlinear disentanglement in natural data with temporal sparse coding. In *International Conference on Learning Representations*, 2021. URL https://openreview.net/forum?id=EbIDjBynYJ8. 2, 7
- [37] Adam Kosiorek, Hyunjik Kim, Yee Whye Teh, and Ingmar Posner. Sequential attend, infer, repeat: Generative modelling of moving objects. In *NIPS*, 2018. 3
- [38] Brenden M Lake, Tomer D Ullman, Joshua B Tenenbaum, and Samuel J Gershman. Building machines that learn and think like people. *Behav. Brain Sci.*, 40, 2017.
- [39] Rosanne Liu, Joel Lehman, Piero Molino, Felipe Petroski Such, Eric Frank, Alex Sergeev, and Jason Yosinski. An intriguing failing of convolutional neural networks and the coordconv solution, 2018. 4, 8, 15
- [40] Francesco Locatello, Ben Poole, Gunnar Rätsch, Bernhard Schölkopf, Olivier Bachem, and Michael Tschannen. Weakly-supervised disentanglement without compromises. arXiv preprint arXiv:2002.02886, 2020. 2
- [41] Francesco Locatello, Dirk Weissenborn, Thomas Unterthiner, Aravindh Mahendran, Georg Heigold, Jakob Uszkoreit, Alexey Dosovitskiy, and Thomas Kipf. Object-centric learning with slot attention, 2020. 3
- [42] Weili Nie, Tero Karras, Animesh Garg, Shoubhik Debnath, Anjul Patney, Ankit Patel, and Animashree Anandkumar. Semi-supervised StyleGAN for disentanglement learning. In *Proceedings of the 37th International Conference on Machine Learning*, pages 7360–7369, 2020. 6, 7
- [43] Erik Nijkamp, Mitch Hill, Tian Han, Song-Chun Zhu, and Ying Nian Wu. On the anatomy of mcmc-based maximum likelihood learning of energy-based models. arXiv preprint arXiv:1903.12370, 2019.
- [44] Ethan Perez, Florian Strub, Harm De Vries, Vincent Dumoulin, and Aaron Courville. Film: Visual reasoning with a general conditioning layer. In *AAAI*, 2018. 15
- [45] Geoffrey Roeder, Luke Metz, and Diedrik P. Kingma. On linear identifiability of learned representations. arXiv preprint arXiv:2007.00810, 2020. 2
- [46] Michal Rolinek, Dominik Zietlow, and Georg Martius. Variational autoencoders pursue pca directions (by accident). In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 12406–12415, 2019.

- [47] Yunfu Song and Zhijian Ou. Learning neural random fields with inclusive auxiliary generators. *arXiv* preprint arXiv:1806.00271, 2018. 3
- [48] Aleksandar Stanić and Jürgen Schmidhuber. R-sqair: Relational sequential attend, infer, repeat. arXiv:1910.05231, 2019. 3
- [49] Sjoerd van Steenkiste, Michael Chang, Klaus Greff, and Jürgen Schmidhuber. Relational neural expectation maximization: Unsupervised discovery of objects and their interactions. arXiv preprint arXiv:1802.10353, 2018. 3
- [50] Sjoerd van Steenkiste, Karol Kurach, and Sylvain Gelly. A case for object compositionality in deep generative models of images. arXiv preprint arXiv:1810.10340, 2018. 2
- [51] Ramakrishna Vedantam, Ian Fischer, Jonathan Huang, and Kevin Murphy. Generative models of visually grounded imagination. In *ICLR*, 2018. 1
- [52] Rishi Veerapaneni, John D Co-Reyes, Michael Chang, Michael Janner, Chelsea Finn, Jiajun Wu, Joshua B Tenenbaum, and Sergey Levine. Entity abstraction in visual model-based reinforcement learning. In *CoRL*, 2019. 3
- [53] Jianwen Xie, Yang Lu, Song-Chun Zhu, and Yingnian Wu. A theory of generative convnet. In *International Conference on Machine Learning*, pages 2635–2644, 2016. 3
- [54] Roland S Zimmermann, Yash Sharma, Steffen Schneider, Matthias Bethge, and Wieland Brendel. Contrastive learning inverts the data generating process. *arXiv preprint arXiv:2102.08850*, 2021. 2

A.1 Unsupervised Learning of Compositional Energy Concepts Appendix

In this supplement, we provide additional empirical visualizations of our approach in Section A.1.1. Next, we provide details on experimental setup in Section A.1.2. The underlying code of the paper can be found at the paper website.

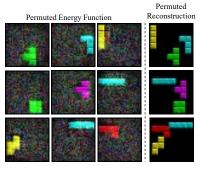


Figure A1: **Tetris Recombination.** Illustration of recombination of energy functions inferred by COMET on the Tetris dataset.

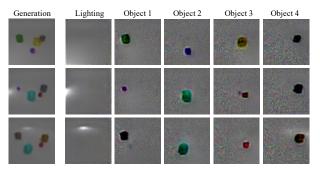


Figure A2: **Light Object Decomposition.** Illustration of decomposing CLEVR Lighting scenes into separate energy functions. COMET is able to decompose a scene into lighting and constituent objects.

Model	$\operatorname{Dim}\left(D\right)$	β	Decoder Dist.	BetaVAE	MIG	MCC
β -VAE (Codebase 1)	64	4	Gaussian	83.57 ± 8.05	10.90 ± 3.80	66.08 ± 2.00
β -VAE (Codebase 2)	64	4	Gaussian	79.99 ± 9.65	7.45 ± 4.58	61.03 ± 5.49
β-VAE (Codebase 1)	32	4	Gaussian	79.77 ± 10.95	7.14 ± 5.44	57.48 ± 6.04
β -VAE (Codebase 2)	32	4	Gaussian	88.01 ± 8.18	12.06 ± 6.05	63.42 ± 5.94
β-VAE (Codebase 1)	256	4	Gaussian	80.76 ± 4.55	10.94 ± 0.58	66.14 ± 1.81
β -VAE (Codebase 2)	256	4	Gaussian	83.93 ± 6.42	7.79 ± 2.41	61.89 ± 2.78
β -VAE (Codebase 1)	64	16	Gaussian	74.71 ± 1.57	9.33 ± 3.72	57.28 ± 2.37
β -VAE (Codebase 2)	64	16	Gaussian	71.30 ± 4.24	6.28 ± 1.18	55.14 ± 3.10
β -VAE (Codebase 1)	64	1	Gaussian	81.61 ± 6.75	6.51 ± 3.38	58.73 ± 6.31
β -VAE (Codebase 2)	64	1	Gaussian	86.42 ± 3.33	9.88 ± 3.27	62.79 ± 3.21
β-VAE (Codebase 1)	64	4	Bernoulli	84.23 ± 3.51	8.96 ± 3.53	61.57 ± 4.09
β -VAE (Codebase 2)	64	4	Bernoulli	93.86 ± 1.74	14.26 ± 2.26	68.34 ± 1.94

Table 3: **Disentanglement Evaluation.** Mean and standard deviation (s.d.) metric scores across 3 random seeds on the Falcor3D dataset. COMET enables better disentanglement across 3 common disentanglement metrics across different runs and seeds for training β -VAE. Note that * denotes that PCA was used as a postprocessing step.

A.1.1 Additional Empirical Results

Qualitative Result. We provide additional empirical visualizations of the results presented in the main paper section 5.2. We illustrate the permuted energy functions in Tetris in Figure A1, and find that we are able to successfully permute different Tetris blocks across different images. We further show that our model is able to decompose an image into both objects and lighting factors in the CLEVR Lighting dataset in Figure A2.

Quantitative Evaluation. We provide an additional comparison of COMET to the β -VAE utilizing a separate codebase across 3 separate seeds in Table 3. We find that across different implementations COMET improves performance over the β -VAE.

A.1.2 Model and Experimental Details

Dataset Details. For the CLEVR dataset, we utilize the dataset generation code from [30] to render images of scenes with 4 objects. For the CLEVR lighting dataset, we also utilize the code

from [30] to render the dataset but increase the lighting jitter to 10.0 for all settings. Finally, for the CLEVR toy dataset, we utilize the default CLEVR dataset generation code but replace the blend files of "sphere", "cylinder", and "cube" with that of "bowling pin", "boot", "toy" and "truck".

Architecture Details. In COMET we utilize a residual network to parameterize an underlying energy function. We illustrate the underlying architecture of the energy function in Figure 4. The energy function takes as input an image at 64×64 resolution and processes the image through a series of residual blocks combined with average pooling to obtain a final output energy. To condition the energy function on an input latent z, we linearly map z to a separate per channel gain and bias in each residual block of the energy network, and use the resultant gain and bias vectors to modulate input features [44]. We remove normalization layers from our residual network.

To infer global factors from an input image, we utilize a convolutional encoder in Figure 5. The convolutional encoder maps an input image through a series of residual convolutional layers to obtain a set of distinct latent vectors. These resultant latent vectors are utilized to condition each separate energy function and correspond to individual global factors.

To infer local object factors in a scene, we utilize a convolutional encoder in combination with a recurrent network with spatial attention. We concatenate a positional embedding to the underlying input image of scene [39]. We then utilize a series of 3 residual layers to downsample input images at 64×64 resolution to a lower resolution 8×8 feature grid. We utilize a LSTM with an attention mechanism [2] to iteratively gather information from this feature grid to obtain a set of object latents representing the scene. We illustrate the overall architecture in Figure 6.

Training Details. Models for each dataset are trained for 12 hours on a single 32GB Volta machine. Models are trained utilizing the Adam optimizer [35] with a learning rate of 3e-4. We utilize 10 steps of optimization to approximate the minimal energy state of an energy function, with each gradient descent step utilizing a scalar multiplier of 1000. When training each model, the gradients are clipped to a magnitude of 1. Images are fit at 64×64 resolution, with a training batch size of 32. We utilize a latent dimension of 16 per component when extracting local factors of variation and a latent dimension of 64 when extracting global factors of variation. For global factor recombination, to obtain high-resolution results on real datasets, we utilize LPIPS loss as a replacement to MSE loss.

3x3 Conv2d, 64
ResBlock Down 64
ResBlock Down 64
ResBlock Down 64
Global Mean Pooling
$\boxed{ \text{Dense} \rightarrow 1 }$

Table 4: Architecture of energy function.

3x3 Conv2d, 64
ResBlock Down 64
ResBlock Down 64
ResBlock Down 64
Global Mean Pooling
${}$ Dense $\rightarrow 64$
${}$ Dense \rightarrow Latents

Table 5: The model architecture used for the convolutional encoder.

3x3 Conv2d, 64
ResBlock Down 64
ResBlock Down 64
ResBlock Down 64
LSTM
Dense \rightarrow 64
$Dense \rightarrow Latents$

Table 6: The model architecture used for the recurrent encoder used in Section 5.2 of the main paper. We utilize a LSTM which operates on the spatial output of the residual network through attention [2].