# DisUnknown: Distilling Unknown Factors for Disentanglement Learning

Sitao Xiang[1,2], Yuming Gu[1,2], Pengda Xiang[1,2], Menglei Chai[3], Hao Li[1,2], Yajie Zhao[2], Mingming He[*2]
[1]University of Southern California, [2]USC Institute for Creative Technologies, [3]Snap Inc.
`sitaoxia@usc.edu`, `{ygu, pxiang}@ict.usc.edu`, `mchai@snap.com`, `hao@hao-li.com`,
`{zhao, he}@ict.usc.edu`

## Abstract

*Disentangling data into interpretable and independent factors is critical for controllable generation tasks. With the availability of labeled data, supervision can help enforce the separation of specific factors as expected. However, it is often expensive or even impossible to label every single factor to achieve fully-supervised disentanglement. In this paper, we adopt a general setting where all factors that are hard to label or identify are encapsulated as a single unknown factor. Under this setting, we propose a flexible weakly-supervised multi-factor disentanglement framework **DisUnknown**, which **Dis**tills **Unknown** factors for enabling multi-conditional generation regarding both labeled and unknown factors. Specifically, a two-stage training approach is adopted to first disentangle the unknown factor with an effective and robust training method, and then train the final generator with the proper disentanglement of all labeled factors utilizing the unknown distillation. To demonstrate the generalization capacity and scalability of our method, we evaluate it on multiple benchmark datasets qualitatively and quantitatively and further apply it to various real-world applications on complicated datasets*

## 1. Introduction

Disentanglement learning is the task of breaking down the tangled high-dimensional data variation into interpretable factors. In the desired disentangled representation, each dimension corresponds to a distinct factor of variables, such that when one factor changes, the others remain unaffected [4]. Disentanglement learning thus enables various downstream tasks such as transfer learning and few-shot learning, as well as challenging controllable image synthesis applications (*e.g.* [50, 16]).

With the availability of fully-labeled data, *supervised disentanglement* has seen much progress [31, 40, 17, 1, 16]. However, ground-truth labels are not always accessible,

while even human labeling could be prohibitively expensive or inconsistent. Thus, fully-supervised approaches often have a hard time generalizing to common scenarios where labels are only partially available or even entirely missing. In light of this, *unsupervised disentanglement* approaches [12, 22, 29, 53, 44] have been proposed to address these challenges. However, most of them rely on the strong assumption that the target data is well-structured enough to be cleanly decoupled into explanatory and recoverable factors. And more importantly, there is no guarantee that these factors could be explicitly controlled with respect to the true intended semantics in specific manipulation scenarios. Therefore, *weakly-supervised disentanglement*, a nice mix of the best of both worlds, has recently become popular for more flexible learning [31, 47, 10, 19]. Unfortunately, although state-of-the-art performance is achieved on certain two-factor class-content disentanglement tasks [10, 19], most existing methods in this category are still unable to extract factor-aware latent representation, which is essential for manipulating individual factors especially when multiple ones are presented. In conclusion, no solution seems completely satisfactory yet on multi-factor disentanglement, due to the limited generalizability and insufficient performance.

In this paper, we propose a weakly-supervised multi-factor disentanglement learning framework, which handles arbitrary numbers of factors through explicit and near-orthogonal latent representation. Given that challenging factors that are hard to label or interpret exist in most tasks, the *key idea* to our approach is a general setting of $N$-factor disentanglement with $N - 1$ factors labeled and a single factor unknown, where all the remaining task-irrelevant or difficult-to-label factors are flexibly encapsulated as one unknown factor. We find such a setting highly effective and practical in real scenarios. Take face motion retargeting as an example, facial expression could be a good candidate for the unknown factor since it is much more difficult to precisely label than others such as the identity and the pose. Thanks to its flexibility, our method naturally adapts to various tasks with varying domains (*e.g.* cartoon and real pho-

---

[*]Corresponding author.
 Project website: https://stormraiser.github.io/disunknown/

tos), data types (*e.g.* images, skeletons, and landmarks), integrity (well-structured or in-the-wild), and label continuity (discrete or continuous).

To this end, our framework consists of two major stages: 1) *Unknown Factor Distillation* and 2) *Multi-Conditional Generation*. Specifically, we extract the unknown factor using an adversarial training method in the first stage, and then embed all labeled factors to the latent space as the second stage, which are used to condition the final generation. The core of our method lies in the joint adversarial training of factor encoders and discriminative classifiers, which explicitly disentangles unknown and known factors without introducing leakage between their disentangled representations.

The performance of our approach is extensively evaluated on several benchmark datasets, both qualitatively and quantitatively. Furthermore, we demonstrate the generalization capacity and practical robustness of the framework on multiple challenging tasks using complicated real-world datasets without any additional manual labeling effort.

Our contributions are: 1) A flexible weakly-supervised disentanglement learning framework that models data as a combination of labeled/unlabeled factors, which scales well to different datasets and benefits various challenging tasks; 2) A two-stage training architecture that explicitly learns disentangled representations for both labeled and unknown semantic factors, enabling mutual exclusive manipulation in the dimension of each factor; 3) A set of learning strategies to improve the effectiveness and robustness of adversarial training throughout our pipeline, which could potentially inspire future research; 4) State-of-the-art performance and wide range of practical uses on multiple challenging tasks including controllable image generation.

## 2. Related Work

**Unsupervised Disentanglement** has become the research focus because it does not require the access to the factors of variation. The pioneering work of InfoGAN [12], an information-theoretic extension to the Generative Adversarial Network framework [21], learns disentangled representations by maximizing the mutual information between the observations and a subset of latents. Considering its training instability and reduced diversity, the Variational Autoencoder (VAE)-based methods [22, 11, 32, 37, 29] are proposed for better performance and reconstruction quality by enforcing a factorized aggregated posterior on the latent space. However, these models are built on the assumption that the observations are independent and identically distributed in the datasets, thus successfully disentangled models may not be identified without any supervision [36]. Some task-specific unsupervised approaches disentangle two or more factors and achieve impressive results, such as image-to-image translation [23, 34, 45] and motion retargeting [52, 63]. These methods do learn disentangled rep-

resentations, relying on specific categories [56, 51, 38, 63], clearly defined domains [23, 34, 45], or well-structured datasets with certain categories [53, 35]. In contrast, our method proposes a general framework, adapting to various tasks, domains, modalities and factor numbers.

**Supervised Disentanglement** requires strong supervision on specific factors of the data. These methods train a subset of the representations to match the known labels using supervised learning [46, 61]. With observed class labels only available for partial data, [24] and [42] propose semi-supervised VAE methods that learn disentangled representation. These supervised methods require large amounts of supervised data that would be expensive to acquire in practice. Although some methods can use synthetic data or data priors to provide full supervision [1, 16, 55], they are limited to processing domain-specific data such as human faces/bodies/hairstyles. Comparing to most supervised methods that only apply to specific tasks, what we propose is a general approach that applies to various applications.

**Weakly-Supervised Disentanglement** has been recently studied to build robust disentangled representations without requiring large amounts of data. Such weak supervision is provided as either known relations between the factors in different samples or ground truth labels of a subset of factors. To avoid explicitly labeling, some methods consider guiding disentanglement by matching pairs of data that share the same underlying factor [47, 31, 24, 5, 10]. By observing a subset of the ground truth factors, some methods perform distribution matching over data and observed factors and supervision is leveraged in style-content disentanglement with available labels for style only [30, 62, 28, 19]. Some of these methods may achieve state-of-the-art performance on certain class-content disentanglement tasks [10, 19], but they cannot ensure factor-aware latent representations for manipulating individual factors. The similar idea of a unified representation of labeled/unlabeled factors has emerged [18]. But we present a general disentanglement learning framework, which benefits various tasks.

## 3. Method

We propose a generic framework for weakly-supervised disentanglement learning and conditional generation. Instead of jointly training the whole system altogether, we take a two-stage approach. In the first stage, excluding all labeled factors, an encoder is trained to extract disentangled representation of the unknown factor from the input data. And in the second stage, with the unknown factor distilled, a conditional generative adversarial network is trained to embed the labeled data into the latent space, which allows independent control over each factor. By isolating the unknown factor from the labeled ones first, this two-stage training helps reduce the overall complexity of the task and improve
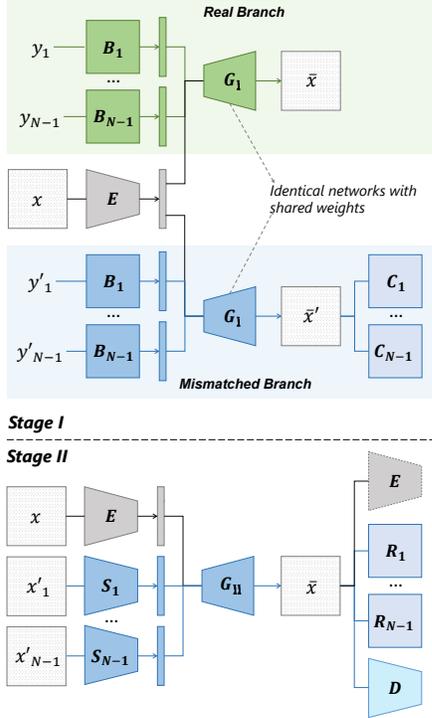
Figure 1: Illustration of our two-stage training architecture.

the effectiveness of labeled factor disentanglement, as will be elaborated in the Training Strategy part in Stage II.

We note that Stage II is fully-supervised, in which missing labels for the unknown factor is provided by Stage I. Thus our method trivially covers the case where all factors are labeled, by dropping Stage I and using only Stage II.

### 3.1. Stage I: Unknown Factor Distillation

This stage trains an *unknown encoder* $E$ that encodes the unknown factor completely and exclusively. It has two parallel branches in Figure 1 (*Stage I*), taking ground truth labels (in the *real branch*) and random labels (in the *mismatched branch*) of all known factors as input, respectively.

Specifically, let there be $N$ factors, with the first $N-1$ ones labeled and the last one unlabeled. $x$ is the training sample, $y = \{y_1, \ldots, y_{N-1}\}$ are the associated *ground truth labels* and $y' = \{y'_1, \ldots, y'_{N-1}\}$ are *random labels* chosen independently of $x$. $E$ is the aforementioned *unknown encoder*, $B = \{B_1, \ldots, B_{N-1}\}$ is a set of *label embedders*, both output normal distributions as in a VAE. $G_I$ is the *Stage-I generator* that generates a sample $\bar{x}$ or $\bar{x}'$ for the real or mismatched branch, respectively, conditioned on $E$ and $B$. $C = \{C_1, \ldots, C_{N-1}\}$ is a set of *classifiers* that predicts the probability distribution of each factor from a generated sample. Both branches share network structures and weights. The loss functions of the two branches are summed. For now, we assume discrete labels, and discuss continuous-valued factors in the supplementary material.

*Real branch:* $B$ map the ground truth labels $y$ to normal distributions. We sample codes from these distributions and feed them to $G_I$, together with the distilled unknown factor from $E$, to generate the reconstructed sample $\bar{x}$.

*Mismatched branch:* By replacing the ground truth labels with random ones $y'_i$, $G_I$ is asked to generate a mixed sample $\bar{x}'$. $C_i$ predicts the ground truth label from the mixed sample, which indicates if any label information is leaked through $E$, since only $E$ has the access to the ground truth factors in $x$. $C$ are implemented as a single multi-class classifier that only branches at the last layer, and are trained with $E$ in an adversarial manner.

**Motivation.** 1) In the real branch, by enforcing a reconstruction loss between the generated sample $\bar{x}$ and the original one $x$, $E$ should include all information not covered by any labeled factor; 2) In the mismatched branch, by minimizing the accuracy of the classifiers $C$ that are trying to predict the ground truth labels from the generated mixed sample $\bar{x}'$, $E$ should exclude any information associated with the labeled factors to avoid label leaking.

**Training Strategy.** As a common problem of adversarial methods, jointly training the adversarial pair of $E$ and $C$ could be unstable. To improve the training robustness, we operate $C$ on samples generated by $G_I$ instead of codes sampled from the distributions produced by $E$ (similar to [14]). This is because, without proper constraints, the distributions in the code space can fluctuate a lot in attempting to prevent the code from being classified. In contrast, with the reconstruction loss in the sample space, the distributions of the generated samples are close to the real ones, which avoids this kind of fluctuation.

As usual, the classifier $C$ minimizes the *negative log-likelihood* (NLL). Let $p$ be a vector representing the probability distribution for a particular factor and $k$ be a class label whose probability is $p_{(k)}$, NLL is defined as:

$$\mathsf{NLL}(p, k) = -\ln p_{(k)}. \tag{1}$$

As the adversarial counterpart, the most obvious choice for the adversarial loss of $E$ is to maximize the NLL loss. However, since NLL is not bounded when the probability $p_{(k)}$ is close to zero, $E$ may prefer to focus on scoring very large NLL values on only a few samples rather than to make every output code equally unclassifiable. Therefore, instead of maximizing the NLL loss, we propose to minimize the *weighted negative log-unlikelihood loss* (NLU):

$$\mathsf{NLU}_q(p, k) = -\frac{1 - q_{(k)}}{q_{(k)}} \ln(1 - p_{(k)}), \tag{2}$$

where $q$ are the reference distributions, which are always taken to be the actual class distributions in the training set for our purpose. In the supplementary material, we show how this definition of NLU loss is derived from the desired

properties that it should be bounded, yield larger gradients on samples farther from equilibrium, and have the same equilibrium point as maximizing the NLL loss.

**Full Objective.** The full training objective on a single sample for Stage I is formulated as:

$$(\mu, \sigma^2) = E(x), \quad e \sim \mathcal{N}(\mu, \text{diag}(\sigma^2)), \tag{3a}$$

$$(\alpha_i, \beta_i^2) = B_i(y_i), \quad b_i \sim \mathcal{N}(\alpha_i, \text{diag}(\beta_i^2)), \tag{3b}$$

$$(\alpha_i', (\beta_i')^2) = B_i(y_i'), \quad b_i' \sim \mathcal{N}(\alpha_i', \text{diag}((\beta_i')^2)), \tag{3c}$$

$$\overline{x} = G_{\text{I}}(e, b_1, \ldots, b_{N-1}), \tag{3d}$$

$$\overline{x}' = G_{\text{I}}(e, b_1', \ldots, b_{N-1}'), \quad p_i = C_i(e, \overline{x}'), \tag{3e}$$

$$\mathcal{L}_C = \sum_i \text{NLL}(p_i, y_i), \tag{3f}$$

$$\mathcal{L}_{GEB} = \text{Rec}(x, \overline{x}) + \lambda_{\text{adv1}} \sum_i \text{NLU}_q(p_i, y_i)$$
$$+ \lambda_{\text{KL}} D_{\text{KL}}(\mathcal{N}(\mu, \text{diag}(\sigma)) || \mathcal{N}(\mathbf{0}, I)) \tag{3g}$$
$$+ \lambda_{\text{KL}} \sum_i D_{\text{KL}}(\mathcal{N}(\alpha_i, \text{diag}(\beta_i^2)) || \mathcal{N}(\mathbf{0}, I))).$$

The square on the variance vectors $\sigma^2$, $\beta_i^2$ and $(\beta_i')^2$ are per-element. $\text{Rec}(x, \overline{x})$ is the reconstruction loss function, which is the mean squared error $||x - \overline{x}||^2$ in our experiments. $D_{KL}$ is the KL-divergence. $C$ are trained in the mismatched branch to minimize $\mathcal{L}_C$, averaged over all samples. $E$, $B$, and $G_{\text{I}}$ jointly minimize $\mathcal{L}_{GEB}$.

## 3.2. Stage II: Multi-Conditional Generation

With the unknown factor distilled in Stage I, this second stage trains encoders $S$ for labeled factors to extract the disentangled representations from the input samples. The final multi-conditional generator $G_\Pi$ accepts conditions for both labeled and unknown factors, and ensures that varying one factor would not affect others in the generated output.

In this stage, as shown in Figure 1 (Stage II), the conditions of the unknown and labeled factors come from training samples $x$ and $\{x_1', \ldots, x_{N-1}'\}$ respectively, all chosen independently. Each $S_i$ of the *labeled-factor encoders* $S = \{S_1, \ldots, S_{N-1}\}$ computes the code for labeled factor $i$ from $x_i'$, while the *unknown encoder* $E$, pre-trained in Stage I, computes the unknown factor code from $x$. The *Stage-II generator* $G_\Pi$ generates a sample $\overline{x}$ conditioned on all the codes (Eq. 5c). On $\overline{x}$, a set of *discriminative classifiers* $R = \{R_1, \ldots, R_{N-1}\}$ are trained to enforce the independent controllability of the labeled factor codes, and the pre-trained $E$ is adopted to ensure the consistency of the unknown factor. In addition, a *discriminator* $D$ is applied to ensure the realism of generated samples, as in GAN.

**Motivation.** Trained on random combinations of input samples, the generator $G_\Pi$ is asked to synthesis a new sample with each factor conditioned by encodings from independent sources. Each classifiers $R_i$ enforces that factor $i$ of $\overline{x}$ is completely and solely controlled by $x_i'$, and by choosing each $x_i'$ randomly and independently we ensure that $S_i$ is the only encoder that can consistently compute factor $i$ of

$x_i'$. The discriminator $D$ makes the distribution of generated samples and real data indistinguishable globally.

**Training Strategy.** Most previous class-conditional GANs differ on how the generated sample is treated by the classifiers. Their classifiers are trained to correctly label the generated sample [43] or to be uncertain about the task [54]. But we go the opposite way: in addition to the NLL loss (Eq. 5e) for classifying the training sample $x$ to the correct labels, our discriminative classifiers $R$ are specifically trained to *not* classify the generated sample $\overline{x}$ correctly, by adding the *unweighted* NLU loss:

$$\text{NLU}(p, k) = -\ln(1 - p_{(k)}). \tag{4}$$

Its rationale is that a conventional classifier oblivious to the generated samples tends to only learn just enough to distinguish one class from the others, which is insufficient to define the full characteristics of that class. However, if we ask the classifier to identify a generated sample as being in the wrong class, in order to tell real and generated samples apart it would be encouraged to gain a more complete understanding of each class.

$G_\Pi$ and $S$ are jointly trained to ensure that the generated sample $\overline{x}$ is classified to the same labels as the inputs $\{x_1', \ldots, x_{N-1}'\}$ (the NLL term in Eq. 5g).

Meanwhile, to enforce that the unlabeled factor is consistently controlled by the code from $E$, we minimize the distance between the encodings of the generated sample $\overline{x}$ and the input $x$, using the fixed $E$ (square error term in Eq. 5g). This further explains why $E$ must be trained in a separate stage from the rest of the system: $E$ is used both for providing the input to the generator and for re-encoding the output to compare against the input. If $E$ is allowed to be updated while this distance is being minimized, it could collapse to a state where it encodes everything to a zero vector.

As for the discriminator $D$, we use LSGAN loss functions [39] (Eq. 5f and the D term in Eq. 5g).

**Full Objective.** Similar to Stage I, the full training objective on a single sample for Stage II is formulated as:

$$(\mu, \sigma^2) = E(x), \quad e \sim \mathcal{N}(\mu, \text{diag}(\sigma^2)), \tag{5a}$$

$$(\alpha_i', (\beta_i')^2) = S_i(x_i'), \quad s_i' \sim \mathcal{N}(\alpha_i', \text{diag}((\beta_i')^2)), \tag{5b}$$

$$\overline{x} = G_\Pi(e, s_1', \ldots, s_{N-1}'), \quad (\overline{\mu}, \overline{\sigma}^2) = E(\overline{x}), \tag{5c}$$

$$p_i = R_i(x), \quad p_i' = R_i(\overline{x}), \tag{5d}$$

$$\mathcal{L}_R = \sum_i (\text{NLL}(p_i, y_i) + \text{NLU}(p_i', y_i')), \tag{5e}$$

$$\mathcal{L}_D = (D(x) - 1)^2 + (D(\overline{x}) + 1)^2, \tag{5f}$$

$$\mathcal{L}_{GS} = ||\overline{\mu} - \mu||^2$$
$$+ \lambda_{\text{adv2}}(D(\overline{x})^2 + \sum_i \text{NLL}(p_i', y_i')) \tag{5g}$$
$$+ \lambda_{\text{KL}} \sum_i D_{\text{KL}}(\mathcal{N}(\alpha_i', \text{diag}((\beta_i')^2)) || \mathcal{N}(\mathbf{0}, I)).$$

Note that while a total of $N$ input samples are required to generate one sample, in practice this can be efficiently done

by computing all factor codes for a whole batch and combining them randomly for generation. Classification labels are permuted accordingly. The classifiers $R$ minimize $\mathcal{L}_R$, the discriminator $D$ minimizes $\mathcal{L}_D$, and the generator $G$ and encoders $S$ jointly minimize $\mathcal{L}_{GS}$.

## 3.3. Implementation Details

For maximum generality we do not favor any specific network architecture. In all our experiments, encoders and generators consist of 3, 4, or 5 stride-2 convolutions for datasets with image sizes of 28, 64, or 128, respectively, followed by 3 fully-connected layers. Discriminators and classifiers have the same convolutional layers but only one fully-connected layer. The convolution feature map depth starts from 32 and doubles after each convolution but does not exceed 256. Fully-connected layers have 512 features.

## 4. Experiments

### 4.1. Datasets and Metrics

**Datasets.** We conduct evaluation experiments on four benchmark datasets: *MNIST* [33], *Fashion-MNIST* (*F-MNIST*) [60], *3D Chairs* [2], and *3D Shapes* [6]. For *MNIST* and *F-MNIST*, we use the standard training/testing split. For *3D Chairs* and *3D Shapes*, we randomly hold out $10\%$ of all images for testing and use the rest for training. In *MNIST* and *F-MNIST*, we take *class* as the labeled factor since only it has labels available. In *3D Chairs* which contains three factors, i.e. *model*, *elevation*, and *azimuth*, we combine *elevation* and *azimuth* in to a single unknown factor of *rotation*. In *3D Shapes* which is fully defined by six labeled factors, i.e. *floor hue*, *wall hue*, *object hue*, *scale*, *shape*, and *orientation*, we select one or more factors as labeled and merge the remaining ones into the unknown factor to train various models for our empirical study.

**Metrics.** We evaluate the disentanglement performance by computing the Mutual Information Gap (MIG) [11] of the encoders. Since factors may contain more than one dimension, the mutual information of each factor is defined as the largest one over all dimensions. Then the MIG is computed as the gap of mutual information between the top two factors. Higher MIGs indicate better disentanglement quality.

### 4.2. Empirical Study

We empirically study how unknown distillation contributes to the disentanglement of labeled factors and enables control over the unknown factor.

**Necessity of the Unknown Factor.** Without the unknown distillation, there is no guarantee that the features represented by the unknown factor remain fixed when altering any labeled ones. To compare, we modify Stage II by replacing the unknown factor code encoded by $E$ with Gaussian noise and removing the feature matching loss $||\overline{\mu} - \mu||^2$

Table 1: Unknown consistency ratios on *3D Shapes* with different unknown factors, w/ and w/o distillation.

| Unknown Factor | w/ Distillation | w/o Distillation |
|---|---|---|
| *Floor hue* | 100.00% | 63.42% |
| *Wall hue* | 100.00% | 55.63% |
| *Object hue* | 100.00% | 68.76% |

Table 2: Labeled consistency ratios and MIG scores on *3D Shapes* with the unknown factor merged from varying numbers of factors. Zero unknown means fully-supervised.

| # Unknown | Ratio | MIG ↑ |
|---|---|---|
| 0 | 100.00% | 0.9501 |
| 1 | 100.00% | 0.9555 |
| 2 | 100.00% | 0.9733 |
| 3 | 100.00% | 0.9718 |
| 4 | 100.00% | 0.9393 |
| 5 | 100.00% | 0.9868 |

Table 3: Mean squared error (MSE) and MIG scores on *3D Shapes* with different unknown factor.

| Unknown Factor | MSE ↓ | MIG ↑ |
|---|---|---|
| *Floor hue* | 0.00049 | 0.9607 |
| *Wall hue* | 0.00063 | 0.9825 |
| *Object hue* | 0.00074 | 0.9766 |
| *Scale* | 0.00062 | 0.9411 |
| *Shape* | 0.00064 | 0.9637 |
| *Orientation* | 0.00064 | 0.9537 |

(Eq. 5g), and train three models on *3D Shapes*, with each selecting *floor hue*, *wall hue*, and *object hue* as the unknown factor, respectively. We generate images using the same random code for the unknown factor and independently-sampled random codes for all labeled factors, and then calculate the ratio of results sharing the same unknown feature, namely *consistency ratio*. Due to the simplicity of *3D Shapes*, these three features can be reliably computed by taking the colors at fixed pixel coordinates. Two colors are considered the same if their $L2$ RGB distance is less than half of the mean distance between two adjacent hue samples in the dataset. We generate 10,000 images for each network, and show the results in Table 1. As can be seen, all ratios reach $100\%$ with distillation, meaning the unknown factor remains unchanged for all test samples. Note that MIGs are not measured here because the disentanglement performance among labeled factors is generally not affected.

**Scope of the Unknown Factor.** In our setting, if there is more than one unknown factor, all these factors will be treated as a whole without individual controllability. However, we can still ensure that the unknown factors are isolated from the labeled ones, and the disentanglement performance of the labeled factors will not be influenced. To

(a) MNIST/ class/ style     (b) F-MNIST/ class/ style     (c) 3D Chairs/ model/ rotation     (d) 3D Shapes/ noted/ others
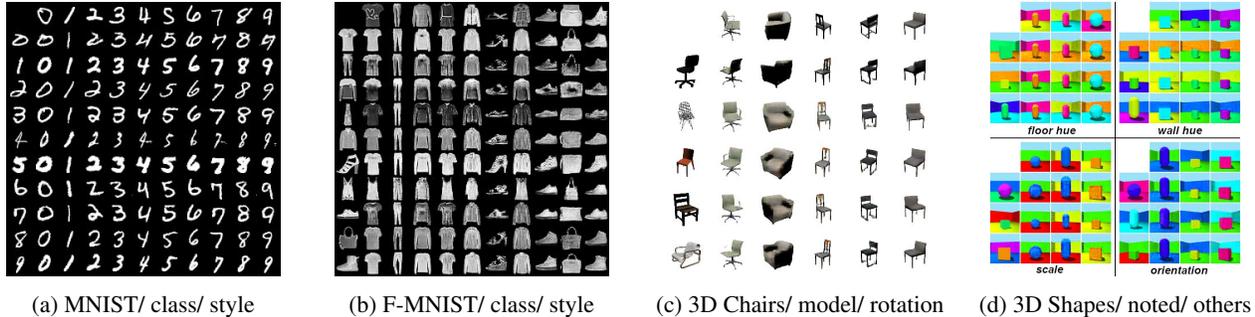
Figure 2: Generated samples on different datasets. The top row and the leftmost column are the input conditions for the labeled and the unknown factors, respectively, annotated as *dataset / labeled / unknown* in the sub-captions.
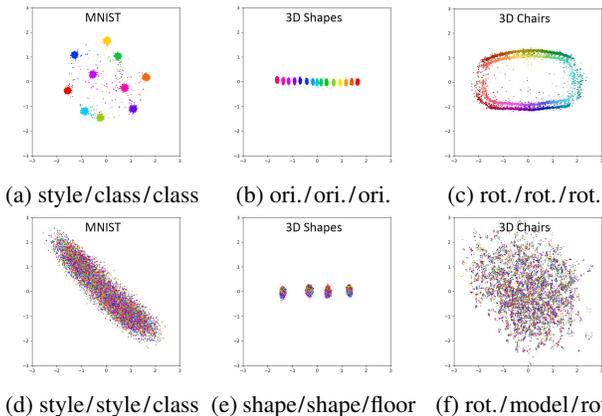


(a) style/class/class     (b) ori./ori./ori.     (c) rot./rot./rot.

(d) style/style/class    (e) shape/shape/floor    (f) rot./model/rot.

Figure 3: Visualizing the disentanglement with test sample distributions. The sub-caption of each figure represents: *unknown factor/ encoding factor/ coloring factor*.

verify this, we train six models on *3D Shapes*: starting all factors labeled, we successively merge *floor hue*, *orientation*, *wall hue*, *scale*, and *shape* into the unknown factor, with *object hue* being the last labeled factor at the end. We measure the consistency ratios as introduced in *Necessity of the Unknown Factor* and MIG scores on *object hue* only in Table 2. Note that all MIG scores are quite close to the upper bound of 1, suggesting good disentanglement quality. **Choice of the Unknown Factor.** We also study the robustness of our method by choosing different factors as the unknown one on *3D Shapes*. The MSE and MIG results, reflecting the consistent performance of reconstruction and disentanglement, respectively, are shown in Table 3.

### 4.3. Results and Visualizations

To demonstrate the quality of our multi-conditional generator, we plot the generated samples with factors controlled by random references on the benchmark datasets. As shown in Figure 2, our method accurately encodes both known (the top row) and unknown (the leftmost column) factors and uses them to independently control the generation.

We also illustrate the disentanglement quality by visu-

alizing the test sample distributions in the code spaces in Figure 3. For each figure, we pick one encoding factor and one coloring factor from all factors, where both factors may or may not be the same. To draw each test sample on the 2D visualization, we generate the 2D position with the encoding factor and the color with the coloring factor. Specifically, we get its factor code using the encoder corresponding to the encoding factor and project it to 2D by selecting two dimensions with the largest variance. Then we draw a point on that 2D projection using the color mapped to its label of the coloring factor. The indication of good disentanglement is that colors should be clearly separated when the encoding and coloring factors are identical, but entirely mixed with no color pattern or bias when they are different.

### 4.4. Comparisons

We compare our approach against the state-of-the-art, including unsupervised [22, 29, 11] and weakly-supervised methods [10, 19]. The weakly-supervised methods are run under the same setting as ours where only one factor is labeled for *MNIST*, *F-MNIST*, and *3D Chairs*. Suggested hyperparameters are used to train these models: $\beta = 4$ for [22]; $\gamma = 10$ on *MNIST* and *F-MNIST*, and $\gamma = 3.2$ on *3D Chairs* for [29]; $\beta = 6$ for [11]; and $\beta = 10$ for [10].

From the results in Table 4, our method achieves substantially higher MIG scores than other methods on all datasets. Since the unsupervised methods [22, 29, 11] are trained without any supervision, comparing with them is somewhat unfair. Nevertheless, this emphasizes the importance of supervision in the disentanglement tasks, which is also reflected by the observation that the weakly-supervised methods consistently outperform the unsupervised ones.

We show a qualitative comparison in Figure 4 which rotates the *3D Chairs* images via traversing the latent code that depicts the azimuth rotation. The unsupervised methods [22, 29, 11] can smoothly change the orientation but fail to preserve the original style (*e.g.* shape, color, etc.). Among the weakly-supervised methods, [10] suffers from over-blurriness, while [19] cannot consistently control the

Table 4: The MIG scores of different disentanglement methods computed on the benchmark datasets.

| Dataset | Unsupervised | | | Weakly-Supervised | | |
|---------|------|------|------|------|------|------|
| | [22] | [29] | [11] | [10] | [19] | Ours |
| MNIST | 0.279 | 0.071 | 0.568 | 0.760 | 0.582 | **0.978** |
| F-MNIST | 0.105 | 0.043 | 0.111 | 0.630 | 0.539 | **0.874** |
| 3D Chairs | 0.031 | 0.098 | 0.115 | 0.212 | 0.284 | **0.404** |



[22]      [29]

[11]      [10]

[19]      Ours

Figure 4: The rotation manipulation comparison on *3D Chairs* by uniformly sampling the latent codes depicting the azimuth rotation. The leftmost column shows the inputs.

orientation. Instead, our method is capable of handling various chair styles and orientations, and achieves better generation quality with the original styles well preserved. Moreover, both weakly-supervised methods are limited to two-factor class-content disentanglement, but our approach is a more flexible multi-factor framework that supports factor-aware latent representation for each individual factor.

## 5. Downstream Tasks

**Portrait Relighting.** We train the network on the dataset combining celebA-HQ [25] and FFHQ [26] by treating the lighting as the labeled factor and the remaining content as unknown. Here, lighting is represented by second-order spherical harmonics coefficients for RGB and estimated with [27, 8]. Figure 5 shows our portrait relighting results.

**Anime Style Transfer.** We train the network on a custom dataset of $106,814$ anime portrait images drawn by $1,139$
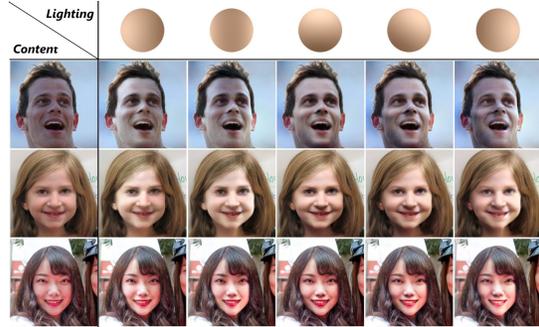


Figure 5: **Portrait relighting.** The top row shows various environment lightings mapped on a sphere. The leftmost column shows input images, and to the right are the re-lit results conditioned by the lightings in the same column.
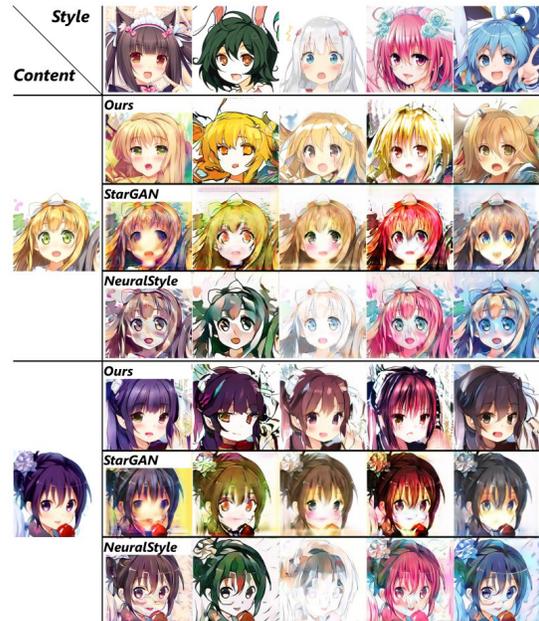


Figure 6: **Anime style transfer.** Each column is conditioned by the example style at the top row. In each group with three rows, the leftmost image is the content and the results are shown to the right. From top to bottom: our method, StarGAN [13], and Neural Style Transfer [20].

artists collected online. The labeled factor is the artists' identity, which is used as the proxy for style. The unlabeled factor is interpreted as the content of the subject. Figure 6 shows our results on transferring style between different anime portrait illustrations, with comparisons to Star-GAN [13] in multi-domain translation and the original Neural Style Transfer [20]. Our method achieves better results with styles more faithful to the examples.

**Landmark-Based Face Reenactment.** We train our disentanglement network on facial landmark coordinates. After the new landmarks are synthesized with our generator, the output face images are translated from the rasterized land-

(a) Fix identity and pose, change facial expression.



(b) Fix identity and facial expression, change pose.

Figure 7: **Face reenactment with expression/pose control**. In each sub-figure, the leftmost column provides the identity and the pose/expression, and the top row provides the expression/pose. The reenactment results are generated with factors conditioned by these inputs.



Figure 8: **Face reenactment with factors from different sources**. The first three rows provide the identity, pose, and expression, respectively. The fourth row shows the results.

marks using the image translation network (*e.g.* [57], [59]). We use FD-GAN in [59] for one-shot image translation. The labeled factors are the identity and the head pose, where the pose is represented by Euler angles, estimated from the landmarks. The unlabeled factor is the facial expression. We train the network on VoxCeleb2 [15]. Figure 7-8 show our face reenactment results with various controls, including editing a single factor (expression/pose) (Figure 7) and mixing all three factors from different sources (Figure 8).
**Skeleton-Based Body Motion Retargeting.** We extract 2D joint coordinates from the driving videos and the actor images. The motion of the driving skeleton and the identity of the actor skeleton are combined to synthesize the target
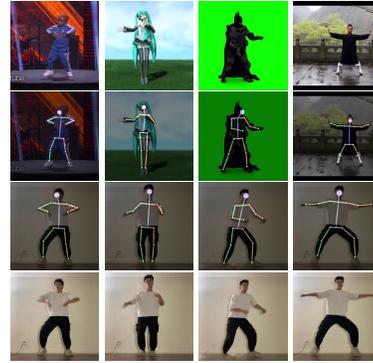


Figure 9: **Body motion retargeting.** From top to bottom in each column: input source frame, extracted source skeleton, transformed skeleton, and generated frame using [48].

skeleton, with motion as the unknown factor. The images are generated using skeleton-guided synthesis (*e.g.* [48], [9]). Figure 9 shows the motion retargeting results on real images trained on Mixamo [41], which demonstrate promising disentanglement between identity and motion.

## 6. Conclusion

We propose *DisUnknown*, a weakly-supervised multi-factor disentanglement learning framework. By distilling unknown factors, it enables independent control over each factor for multi-conditional generation. Our approach achieves state-of-the-art performance compared to existing unsupervised and weakly-supervised methods on multiple benchmark datasets. We further demonstrate its generalization capacity through various downstream tasks. Moreover, as a general framework, it can easily carry over to other modalities (*e.g.* text, audio) and help improve the stability of other tasks with our adversarial training strategies.

# References

[1] Kfir Aberman, Rundi Wu, Dani Lischinski, Baoquan Chen, and Daniel Cohen-Or. Learning character-agnostic motion for motion retargeting in 2d. *ACM Trans. Graph.*, 38(4):75:1–75:14, 2019.

[2] Mathieu Aubry, Daniel Maturana, Alexei A. Efros, Bryan C. Russell, and Josef Sivic. Seeing 3D Chairs: Exemplar Part-Based 2D-3D Alignment Using a Large Dataset of CAD Models. In *CVPR 2014*, pages 3762–3769, 2014.

[3] Tadas Baltrusaitis, Amir Zadeh, Yao Chong Lim, and Louis-Philippe Morency. Openface 2.0: Facial behavior analysis toolkit. In *13th IEEE International Conference on Automatic Face & Gesture Recognition, FG 2018, Xi'an, China, May 15-19, 2018*, pages 59–66. IEEE Computer Society, 2018.

[4] Yoshua Bengio, Aaron C. Courville, and Pascal Vincent. Representation learning: A review and new perspectives. *IEEE Trans. Pattern Anal. Mach. Intell.*, 35(8):1798–1828, 2013.

[5] Diane Bouchacourt, Ryota Tomioka, and Sebastian Nowozin. Multi-level variational autoencoder: Learning disentangled representations from grouped observations. In Sheila A. McIlraith and Kilian Q. Weinberger, editors, *AAAI 2018*, pages 2095–2102, 2018.

[6] Chris Burgess and Hyunjik Kim. 3d shapes dataset. https://github.com/deepmind/3dshapes-dataset/, 2018.

[7] Qiong Cao, Li Shen, Weidi Xie, Omkar M. Parkhi, and Andrew Zisserman. Vggface2: A dataset for recognising faces across pose and age. In *13th IEEE International Conference on Automatic Face & Gesture Recognition, FG 2018, Xi'an, China, May 15-19, 2018*, pages 67–74. IEEE Computer Society, 2018.

[8] Menglei Chai, Linjie Luo, Kalyan Sunkavalli, Nathan Carr, Sunil Hadap, and Kun Zhou. High-quality hair modeling from a single portrait photo. *ACM Trans. Graph.*, 34(6):204:1–204:10, 2015.

[9] Caroline Chan, Shiry Ginosar, Tinghui Zhou, and Alexei A. Efros. Everybody dance now. In *ICCV 2019*, pages 5932–5941, 2019.

[10] Junxiang Chen and Kayhan Batmanghelich. Weakly supervised disentanglement by pairwise similarities. In *AAAI 2020*, pages 3495–3502, 2020.

[11] Tian Qi Chen, Xuechen Li, Roger B. Grosse, and David Duvenaud. Isolating sources of disentanglement in variational autoencoders. In *NeurIPS 2018*, pages 2615–2625, 2018.

[12] Xi Chen, Yan Duan, Rein Houthooft, John Schulman, Ilya Sutskever, and Pieter Abbeel. InfoGAN: Interpretable Representation Learning by Information Maximizing Generative Adversarial Nets. In *NeurIPS 2016*, pages 2172–2180, 2016.

[13] Yunjey Choi, Min-Je Choi, Munyoung Kim, Jung-Woo Ha, Sunghun Kim, and Jaegul Choo. StarGAN: Unified Generative Adversarial Networks for Multi-Domain Image-to-Image Translation. In *CVPR 2018*, pages 8789–8797, 2018.

[14] Ju-Chieh Chou, Cheng-chieh Yeh, Hung-yi Lee, and Lin-Shan Lee. Multi-target voice conversion without parallel data by adversarially learning disentangled audio representations. In *Interspeech 2018*, pages 501–505, 2018.

[15] Joon Son Chung, Arsha Nagrani, and Andrew Zisserman. VoxCeleb2: Deep Speaker Recognition. In *Interspeech 2018*, pages 1086–1090, 2018.

[16] Yu Deng, Jiaolong Yang, Dong Chen, Fang Wen, and Xin Tong. Disentangled and Controllable Face Image Generation via 3D Imitative-Contrastive Learning. In *CVPR 2020*, pages 5153–5162, 2020.

[17] Zunlei Feng, Xinchao Wang, Chenglong Ke, Anxiang Zeng, Dacheng Tao, and Mingli Song. Dual swap disentangling. In *NeurIPS 2018*, pages 5898–5908, 2018.

[18] Zunlei Feng, Zhenyun Yu, Yongcheng Jing, Sai Wu, Mingli Song, Yezhou Yang, and Junxiao Jiang. Interpretable partitioned embedding for intelligent multi-item fashion outfit composition. *ACM Trans. Multim. Comput. Commun. Appl.*, 15(2s):61:1–61:20, 2019.

[19] Aviv Gabbay and Yedid Hoshen. Demystifying inter-class disentanglement. In *ICLR 2020*, 2020.

[20] Leon A. Gatys, Alexander S. Ecker, and Matthias Bethge. Image style transfer using convolutional neural networks. In *CVPR 2016*, pages 2414–2423, 2016.

[21] Ian J. Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron C. Courville, and Yoshua Bengio. Generative adversarial networks. *CoRR*, abs/1406.2661, 2014.

[22] Irina Higgins, Loïc Matthey, Arka Pal, Christopher Burgess, Xavier Glorot, Matthew Botvinick, Shakir Mohamed, and Alexander Lerchner. beta-VAE: Learning Basic Visual Concepts with a Constrained Variational Framework. In *ICLR 2017*, 2017.

[23] Xun Huang, Ming-Yu Liu, Serge J. Belongie, and Jan Kautz. Multimodal unsupervised image-to-image translation. In *ECCV 2018*, volume 11207, pages 179–196, 2018.

[24] Theofanis Karaletsos, Serge J. Belongie, and Gunnar Rätsch. When crowds hold privileges: Bayesian unsupervised representation learning with oracle constraints. In *ICLR 2016*, 2016.

[25] Tero Karras, Timo Aila, Samuli Laine, and Jaakko Lehtinen. Progressive Growing of GANs for Improved Quality, Stability, and Variation. In *ICLR 2018*, 2018.

[26] Tero Karras, Samuli Laine, and Timo Aila. A style-based generator architecture for generative adversarial networks. In *CVPR 2019*, pages 4401–4410, 2019.

[27] Ira Kemelmacher-Shlizerman and Ronen Basri. 3D Face Reconstruction from a Single Image Using a Single Reference Face Shape. *IEEE Trans. Pattern Anal. Mach. Intell.*, 33(2):394–405, 2011.

[28] Bo-Kyeong Kim, Sungjin Park, Geon-min Kim, and Soo-Young Lee. Semi-supervised disentanglement with independent vector variational autoencoders. *CoRR*, abs/2003.06581, 2020.

[29] Hyunjik Kim and Andriy Mnih. Disentangling by factorising. In *ICML 2018*, volume 80, pages 2654–2663, 2018.

[30] Diederik P. Kingma, Shakir Mohamed, Danilo Jimenez Rezende, and Max Welling. Semi-supervised learning with deep generative models. In *NeurIPS 2014*, pages 3581–3589, 2014.

[31] Tejas D. Kulkarni, William F. Whitney, Pushmeet Kohli, and Joshua B. Tenenbaum. Deep convolutional inverse graphics network. In *NeurIPS 2015*, pages 2539–2547, 2015.

[32] Abhishek Kumar, Prasanna Sattigeri, and Avinash Balakrishnan. Variational inference of disentangled latent concepts from unlabeled observations. In *ICLR 2018*, 2018.

[33] Yann LeCun, Léon Bottou, Yoshua Bengio, and Patrick Haffner. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998.

[34] Hsin-Ying Lee, Hung-Yu Tseng, Qi Mao, Jia-Bin Huang, Yu-Ding Lu, Maneesh Singh, and Ming-Hsuan Yang. DRIT++: diverse image-to-image translation via disentangled representations. *Int. J. Comput. Vis.*, 128(10):2402–2417, 2020.

[35] Yuheng Li, Krishna Kumar Singh, Utkarsh Ojha, and Yong Jae Lee. MixNMatch: Multifactor Disentanglement and Encoding for Conditional Image Generation. In *CVPR 2020*, pages 8036–8045, 2020.

[36] Francesco Locatello, Stefan Bauer, Mario Lucic, Gunnar Rätsch, Sylvain Gelly, Bernhard Schölkopf, and Olivier Bachem. Challenging common assumptions in the unsupervised learning of disentangled representations. In *ICML 2019*, volume 97, pages 4114–4124, 2019.

[37] Romain Lopez, Jeffrey Regier, Michael I. Jordan, and Nir Yosef. Information constraints on auto-encoding variational bayes. In *NeurIPS 2018*, pages 6117–6128, 2018.

[38] Dominik Lorenz, Leonard Bereska, Timo Milbich, and Björn Ommer. Unsupervised part-based disentangling of object shape and appearance. In *CVPR 2019*, pages 10955–10964, 2019.

[39] Xudong Mao, Qing Li, Haoran Xie, Raymond Y. K. Lau, Zhen Wang, and Stephen Paul Smolley. Least squares generative adversarial networks. In *ICCV 2017*, pages 2813–2821, 2017.

[40] Michaël Mathieu, Junbo Jake Zhao, Pablo Sprechmann, Aditya Ramesh, and Yann LeCun. Disentangling factors of variation in deep representations using adversarial training. *CoRR*, abs/1611.03383, 2016.

[41] Mixamo. Mixamo. `https://www.mixamo.com/`.

[42] Siddharth Narayanaswamy, Brooks Paige, Jan-Willem van de Meent, Alban Desmaison, Noah D. Goodman, Pushmeet Kohli, Frank D. Wood, and Philip H. S. Torr. Learning disentangled representations with semi-supervised deep generative models. In *NeurIPS 2017*, pages 5925–5935, 2017.

[43] Augustus Odena, Christopher Olah, and Jonathon Shlens. Conditional Image Synthesis with Auxiliary Classifier GANs. In *ICML 2017*, volume 70, pages 2642–2651, 2017.

[44] Ori Press, Tomer Galanti, Sagie Benaim, and Lior Wolf. Emerging disentanglement in auto-encoder based unsupervised image content transfer. In *ICLR 2019*, 2019.

[45] Ori Press, Tomer Galanti, Sagie Benaim, and Lior Wolf. Emerging disentanglement in auto-encoder based unsupervised image content transfer. *CoRR*, abs/2001.05017, 2020.

[46] Scott E. Reed, Kihyuk Sohn, Yuting Zhang, and Honglak Lee. Learning to disentangle factors of variation with manifold interaction. In *ICML 2014*, volume 32, pages 1431–1439, 2014.

[47] Scott E. Reed, Yi Zhang, Yuting Zhang, and Honglak Lee. Deep visual analogy-making. In *NeurIPS 2015*, pages 1252–1260, 2015.

[48] Jian Ren, Menglei Chai, Sergey Tulyakov, Chen Fang, Xiaohui Shen, and Jianchao Yang. Human motion transfer from poses in the wild. In *ECCV 2020 Workshops*, volume 12537, pages 262–279, 2020.

[49] Andreas Rössler, Davide Cozzolino, Luisa Verdoliva, Christian Riess, Justus Thies, and Matthias Nießner. Faceforensics++: Learning to detect manipulated facial images. In *2019 IEEE/CVF International Conference on Computer Vision, ICCV 2019, Seoul, Korea (South), October 27 - November 2, 2019*, pages 1–11. IEEE, 2019.

[50] Yujun Shen, Ceyuan Yang, Xiaoou Tang, and Bolei Zhou. InterFaceGAN: Interpreting the Disentangled Face Representation Learned by GANs. *CoRR*, abs/2005.09635, 2020.

[51] Zhixin Shu, Mihir Sahasrabudhe, Riza Alp Güler, Dimitris Samaras, Nikos Paragios, and Iasonas Kokkinos. Deforming autoencoders: Unsupervised disentangling of shape and appearance. In *ECCV 2018*, volume 11214, pages 664–680, 2018.

[52] Aliaksandr Siarohin, Stéphane Lathuilière, Sergey Tulyakov, Elisa Ricci, and Nicu Sebe. First order motion model for image animation. In *NeurIPS 2019*, pages 7135–7145, 2019.

[53] Krishna Kumar Singh, Utkarsh Ojha, and Yong Jae Lee. FineGAN: Unsupervised Hierarchical Disentanglement for Fine-Grained Object Generation and Discovery. In *CVPR 2019*, pages 6490–6499, 2019.

[54] Jost Tobias Springenberg. Unsupervised and semi-supervised learning with categorical generative adversarial networks. In *ICLR 2016*, 2016.

[55] Zhentao Tan, Menglei Chai, Dongdong Chen, Jing Liao, Qi Chu, Lu Yuan, Sergey Tulyakov, and Nenghai Yu. MichiGAN: multi-input-conditioned hair image generation for portrait editing. *ACM Trans. Graph.*, 39(4):95, 2020.

[56] Luan Tran, Xi Yin, and Xiaoming Liu. Disentangled representation learning GAN for pose-invariant face recognition. In *CVPR 2017*, pages 1283–1292, 2017.

[57] Ting-Chun Wang, Ming-Yu Liu, Jun-Yan Zhu, Andrew Tao, Jan Kautz, and Bryan Catanzaro. High-Resolution Image Synthesis and Semantic Manipulation With Conditional GANs. In *CVPR 2018*, pages 8798–8807, 2018.

[58] Olivia Wiles, A. Sophia Koepke, and Andrew Zisserman. X2face: A network for controlling face generation using images, audio, and pose codes. In Vittorio Ferrari, Martial Hebert, Cristian Sminchisescu, and Yair Weiss, editors, *Computer Vision - ECCV 2018 - 15th European Conference, Munich, Germany, September 8-14, 2018, Proceedings, Part XIII*, volume 11217 of *Lecture Notes in Computer Science*, pages 690–706. Springer, 2018.

[59] Sitao Xiang, Yuming Gu, Pengda Xiang, Mingming He, Koki Nagano, Haiwei Chen, and Hao Li. One-shot identity-preserving portrait reenactment. *CoRR*, abs/2004.12452, 2020.

[60] Han Xiao, Kashif Rasul, and Roland Vollgraf. Fashion-MNIST: a Novel Image Dataset for Benchmarking Machine Learning Algorithms. *CoRR*, abs/1708.07747, 2017.

[61] Taihong Xiao, Jiapeng Hong, and Jinwen Ma. DNA-GAN: learning disentangled representations from multi-attribute images. In *ICLR 2018*, 2018.

[62] Jimei Yang, Scott E. Reed, Ming-Hsuan Yang, and Honglak Lee. Weakly-supervised Disentangling with Recurrent Transformations for 3D View Synthesis. In *NeurIPS 2015*, pages 1099–1107, 2015.

[63] Zhuoqian Yang, Wentao Zhu, Wayne Wu, Chen Qian, Qiang Zhou, Bolei Zhou, and Chen Change Loy. TransMoMo: Invariance-Driven Unsupervised Video Motion Retargeting. In *CVPR 2020*, pages 5305–5314, 2020.

# A. Method Details

## A.1. Derivation of the Negative Log Unlikelihood

Here we show how the form of the adversarial loss of the classifiers are chosen.

First, we consider the assumed equilibrium of the adversarial training. In Stage I, the goal is that the encoder's output should not contain any information about the labeled factors. So, each classifier $C_i$ can at best make a guess, and since there is no way to distinguish between inputs from different classes, it should give the same output class distribution for every sample.

Assume that, as commonly done, the output distribution of the classifier is computed by taking the softmax of a vector $t = (t_{(1)}, t_{(2)}, \ldots, t_{(m)})$, where $m$ is the number of classes for the factor of concern. Let $\mathrm{Sm}$ denote softmax, we have:

$$\frac{\partial}{\partial t_{(i)}} \mathsf{NLL}(\mathrm{Sm}(t), k) \tag{6}$$

$$= \frac{\partial}{\partial t_{(i)}} - \ln \frac{e^{t_{(k)}}}{\sum_j e^{t_{(j)}}}$$

$$= \frac{\partial}{\partial t_{(i)}} (-t_{(k)} + \ln \sum_j e^{t_{(j)}})$$

$$= - \delta_{ik} + \frac{e^{t_{(i)}}}{\sum_j e^{t_{(j)}}}$$

$$= - \delta_{ik} + \mathrm{Sm}(t, i),$$

where $\delta_{ik} = 1$ when $i = k$ and $\delta_{ik} = 0$ otherwise. At equilibrium, the expectation of gradient over the whole dataset should be zero. Let $q = (q_{(1)}, q_{(2)}, \ldots, q_{(m)})$ be the class frequency in the dataset. Then we must have

$$\frac{\partial}{\partial t_{(i)}} \sum_k q_{(k)} \mathsf{NLL}(\mathrm{Sm}(t), k) \tag{7}$$

$$= \sum_k q_{(k)} (-\delta_{ik} + \mathrm{Sm}(t, i))$$

$$= \mathrm{Sm}(t, i) - q_{(i)}$$

$$= 0.$$

That is, $\mathrm{Sm}(t, i) = q_i$. So, at the assumed equilibrium, the classifier should give the class distribution in the dataset as the output for any input.

If the adversarial objective is to maximize the NLL loss of the classifier, then naturally at this assumed equilibrium the expected gradient of the adversarial loss function is also

zero. But, consider the second-order derivatives:

$$\frac{\partial^2}{\partial t_{(i)}^2} \mathsf{NLL}(\mathrm{Sm}(t), k) \tag{8}$$

$$= \frac{\partial}{\partial t_{(i)}} (-\delta_{ik} + \mathrm{Sm}(t, i))$$

$$= \mathrm{Sm}(t, i)(1 - \mathrm{Sm}(t, i))$$

$$> 0,$$

the adversarial loss function has a local minimum with respect to $t$ at the assumed equilibrium, while the objective is to maximize this function. So in the proximity of the assumed equilibrium, the adversarial objective actually pushes the networks away from the equilibrium. Furthermore, consider the L1 norm of the gradient:

$$|| \frac{\partial}{\partial t} \mathsf{NLL}(\mathrm{Sm}(t), k) ||_1 \tag{9}$$

$$= \sum_i | - \delta_{ik} + \mathrm{Sm}(t, i)|$$

$$= \sum_{i \neq k} \mathrm{Sm}(t, i) + (1 - \mathrm{Sm}(t, k))$$

$$= 2 - 2 \cdot \mathrm{Sm}(t, k),$$

the gradient is larger when the NLL loss is larger, and NLL is not bounded above. If the adversarial objective is to maximize the NLL, it can accelerate towards infinity, which causes strong instability.

Remember that a basic trick in vanilla GAN is that instead of letting the generator maximize

$$- \ln(1 - D(G(z))), \tag{10}$$

we let it minimize

$$- \ln(D(G(z))). \tag{11}$$

In a similar vein, instead of maximizing

$$\mathsf{NLL}(p, k) = - \ln p_{(k)}, \tag{12}$$

we can minimize what we call "negative log unlikelihood"

$$\mathsf{NLU}(p, k) = - \ln(1 - p_{(k)}). \tag{13}$$

The derivatives are computed as:

$$\frac{\partial}{\partial t_{(k)}} \mathsf{NLU}(\mathrm{Sm}(t), k) \tag{14}$$

$$= \frac{\partial}{\partial t_{(k)}} - \ln(1 - \frac{e^{t_{(k)}}}{\sum_j e^{t_{(j)}}})$$

$$= - \frac{\sum_j e^{t_{(j)}}}{\sum_j e^{t_{(j)}} - e^{t_{(k)}}} \cdot - \frac{e^{t_{(k)}} (\sum_j e^{t_{(j)}} - e^{t_{(k)}})}{(\sum_j e^{t_{(j)}})^2}$$

$$= \mathrm{Sm}(t, k),$$

$$\frac{\partial}{\partial t_{(i)}}\mathrm{NLU}(\mathrm{Sm}(t),k) \quad (i \neq k) \tag{15}$$

$$=\frac{\partial}{\partial t_{(i)}} - \ln(1 - \frac{e^{t_{(k)}}}{\sum_j e^{t_{(j)}}})$$

$$=-\frac{\sum_j e^{t_{(j)}}}{\sum_j e^{t_{(j)}} - e^{t_{(k)}}} \cdot -\frac{-e^{t_{(i)}}e^{t_{(k)}}}{(\sum_j e^{t_{(j)}})^2}$$

$$=-\frac{\mathrm{Sm}(t,k)\cdot\mathrm{Sm}(t,i)}{1 - \mathrm{Sm}(t,k)}.$$

At the assumed equilibrium $\mathrm{Sm}(t) = q$, where $q$ is the class frequency in the dataset, these evaluate to

$$\left.\frac{\partial}{\partial t_{(k)}}\mathrm{NLU}(\mathrm{Sm}(t),k)\right|_{\mathrm{Sm}(t)=q} \tag{16}$$

$$=q_{(k)},$$

$$\left.\frac{\partial}{\partial t_{(i)}}\mathrm{NLU}(\mathrm{Sm}(t),k)\right|_{\mathrm{Sm}(t)=q} \quad (i \neq k) \tag{17}$$

$$=-\frac{q_{(k)}\cdot q_{(i)}}{1 - q_{(k)}}.$$

If the classes are not evenly distributed in the dataset, this may not satisfy the condition that the assumed equilibrium is a stationary point of $\sum_k q_{(k)}\cdot\mathrm{NLU}(\mathrm{Sm}(t),k)$. To achieve this, we need to properly weight the NLU by class. We can do this by scaling Equation 16 and Equation 17 to match Equation 6. We define the weighted negative log unlikelihood function as:

$$\mathrm{NLU}_q(p,k) = -\frac{1 - q_{(k)}}{q_{(k)}}\ln(1 - p_{(k)}). \tag{18}$$

Then we have, at the assumed equilibrium:

$$\left.\frac{\partial}{\partial t_{(i)}}\sum_k q_{(k)}\cdot\mathrm{NLU}_q(\mathrm{Sm}(t),k)\right|_{\mathrm{Sm}(t)=q} \tag{19}$$

$$=\sum_{k\neq i}-q_{(k)}\cdot\frac{1 - q_{(k)}}{q_{(k)}}\cdot\frac{q_{(k)}\cdot q_{(i)}}{1 - q_{(k)}} + q_{(i)}\cdot\frac{1 - q_{(i)}}{q_{(i)}}\cdot q_{(i)}$$

$$=\sum_{k\neq i}-q_{(k)}q_{(i)} + q_{(i)}(1 - q_{(i)})$$

$$=-(1 - q_{(i)})q_{(i)} + q_{(i)}(1 - q_{(i)})$$

$$=0.$$

And the L1 norm of the gradient would be:

$$\|\frac{\partial}{\partial t}\mathrm{NLU}_q(\mathrm{Sm}(t),k)\|_1 \tag{20}$$

$$=\sum_i \frac{1 - q_{(k)}}{q_{(k)}}\cdot\mathrm{Sm}(t,k)\left(1 + \sum_{i\neq k}\frac{\mathrm{Sm}(t,i)}{1 - \mathrm{Sm}(t,k)}\right)$$

$$=2\cdot\frac{1 - q_{(k)}}{q_{(k)}}\cdot\mathrm{Sm}(t,k),$$

which equals to Equation 9 at the assumed equilibrium, and has the desired property that a smaller value of $\mathrm{Sm}(t,k)$ gives a smaller gradient. Evaluating the second derivative at the assumed equilibrium gives

$$\left.\frac{\partial^2}{\partial t_{(k)}^2}\mathrm{NLU}(\mathrm{Sm}(t),k)\right|_{\mathrm{Sm}(t)=q} \tag{21}$$

$$=q_{(k)}(1 - q_{(k)}),$$

$$\left.\frac{\partial^2}{\partial t_{(i)}^2}\mathrm{NLU}(\mathrm{Sm}(t),k)\right|_{\mathrm{Sm}(t)=q} \quad (i \neq k) \tag{22}$$

$$=\frac{q_{(k)}q_{(i)}(2q_{(i)} + q_{(k)} - 1)}{(1 - q_{(k)})^2},$$

$$\left.\frac{\partial^2}{\partial t_{(i)}^2}\sum_k q_{(k)}\mathrm{NLU}_q(\mathrm{Sm}(t),k)\right|_{\mathrm{Sm}(t)=q} \tag{23}$$

$$=q_{(i)}(1 - q_{(i)})^2 + \sum_{k\neq i}\frac{q_{(k)}q_{(i)}(2q_{(i)} + q_{(k)} - 1)}{1 - q_{(k)}}$$

$$=q_{(i)}\left((1 - q_{(i)})^2 + \sum_{k\neq i}(\frac{2q_{(k)}q_{(i)}}{1 - q_{(k)}} - q_{(k)})\right)$$

$$>q_{(i)}\left((1 - q_{(i)})^2 + \sum_{k\neq i}q_{(k)}(2q_{(i)} - 1)\right)$$

$$=q_{(i)}\left((1 - q_{(i)})^2 + (1 - q_{(i)})(2q_{(i)} - 1)\right)$$

$$=q_{(i)}(1 - q_{(i)})(1 - q_{(i)} + 2q_{(i)} - 1)$$

$$=q_{(i)}^2(1 - q_{(i)})$$

$$>0.$$

So the assumed equilibrium is indeed a local minimum of the adversarial loss function we want to minimize.

### A.2. Sample-Space Classification

The Stage I training procedure of generating samples from random labels for classification is not straightforward to understand, and here we give an explanation.
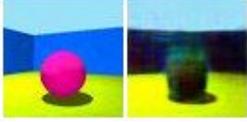
Figure 10: Ill-formed sample generated from mean labeled code.

The distribution of the encoder's output in the code space has few constraints, and there can be different networks that give very different distribution in the code space but are nevertheless essentially equivalent. For example, assume that the last layer of the encoder and the first layer of the generator are linear and the dimension of the code space is $d$. Then we can take an invertible $d \times d$ matrix $M$. We multiply the last layer weights of the encoder by $M$ on the right and multiply the first layer weights of the generator by $M^{-1}$ on the left. In terms of reconstruction, the modified network gives the exact same result as the original, but the code distribution in the code space has been transformed.

This becomes a problem with adversarial training: if the classifier operates in the code space, then to avoid being successfully classified, instead of removing information about the labeled factors from its output, the unknown factor encoder can change its output distribution to confuse the classifier, which results in the code distribution fluctuating constantly in the code space.

In contrast, in the sample space, the distribution of generated samples is anchored to the distribution of training samples and cannot change freely. So, operating the classifier in the sample space can potentially reduce fluctuation in the code space and improve stability. Then the question is which samples should be the input to the classifier.

The reconstructed sample cannot serve as the input to the classifier, since it must contain full information of the input sample, including those about the labeled factors, which is in conflict with the adversarial objective of making the classifier unable to classify by the labeled factor.

Another choice is to combine the unknown code with some kind of "neutral" labeled code, for example, the mean of all label embeddings. The problem is that the "mean" labeled code may not be "typical": as can be seen in Figure 13, in the *3D Shapes* dataset, the network learns that the ten classes of each color attribute are arranged like a circle, but no samples are distributed near the center of the circle. In this case, the mean labeled code does not produce a well-formed sample. An example is shown in Figure 10: on the left is the input. The floor hue is the unknown factor. In the generated sample on the right, all labeled factors have been replaced by the mean embedding, and the generated sample is ill-formed.

The result is that the encoder can encode labeled information about the input sample without being detected:

the generator can easily recognize an invalid mean labeled code, and when it receives one it generates ill-formed samples so that the classifier cannot classify the generated sample, regardless of what the encoder has encoded.

So the unknown code used for generating the input to the classifier must be a typical code but at the same time independent from the input sample. Thus, we use the embedding of a random label chosen independently from the input sample.

### A.3. Continuous-Valued Factor Disentanglement

When presenting our method, we assumed that all labels are discrete, class-type labels. Here we discuss the treatment of continuous-valued factors.

Continuous-valued labels are usually associated with regression problems. So it is reasonable to first attempt to use regressors in place of the classifiers. But there are obvious problems with this approach. Consider training samples $x^{(i)}$ each associated with a single, real-valued label $y^{(i)}$. For each $x^{(i)}$, compute $\overline{x}^{(i)\prime}$ as in Equation 3 in the main text. The regressor $C$ should minimize some kind of distance between $C(\overline{x}^{(i)\prime})$ and $y^{(i)}$, say squared distance $(C(\overline{x}^{(i)\prime}) - y^{(i)})^2$. Then, in the assumed equilibrium where the encoder does not encode any information about the labeled factor, the regressor can only make a guess. To minimize the expected loss, the best guess should be the mean of all $y^{(i)}$. Let $y^* = \sum_{i=1}^{n} y^{(i)}$. Then if there exists training sample $x^*$ whose label is $y^*$ or very close to $y^*$, any adversarial training would not guarantee to prevent the encoder from encoding full information of $x^*$: there is no way to distinguish whether the regressor is giving a $y^*$ because it has detected labeled information in its input, or it is giving a $y^*$ because it detected no such information and is making a guess.

While we have not stated explicitly, we have already provided the solution to working with continuous-valued factors: note that in the *3D Chairs* dataset, the rotation angle is not a true category-type factor, but a quantized continuous-valued factor. Treating it as a category-type factor gives satisfactory results. Similarly, for any continuous-valued factor, we can always divide its range into a suitable number of buckets and quantize the factor into discrete labels. Generally, a few dozen buckets would work fine.

## B. Data and Metrics

### B.1. Data

The image size and list of factors of the datasets are given in Table 5, with the number of possible values of each factor and the length of code we use for the encoder of that factor.

Table 5: Datasets used for evaluation and comparison.

| Dataset | Factor | # of Values | Code Size |
|---|---|---|---|
| MNIST $28 \times 28 \times 1$ | Class (Style) | 10 N/A | 10 64 |
| F-MNIST $28 \times 28 \times 1$ | Class (Style) | 10 N/A | 10 64 |
| 3D Chairs $128 \times 128 \times 3$ | Model Elevation Azimuth | 1393 2 31 | 512 2 2 |
| 3D Shapes $64 \times 64 \times 3$ | Floor hue Wall hue Object hue Scale Shape Orientation | 10 10 10 8 4 15 | 8 8 8 8 8 8 |

## B.2. Metrics

The Mutual Information Gap (MIG) was originally proposed for the unsupervised setting. We made some adjustments to the computation of MIG to suit the weakly-supervised setting and to allow multi-dimensional code spaces for each factor.

Let $N$ be the number of factors, $L_i$ be the (discrete) random variable representing the label of factor $i$ and $X_i$ be the vector-valued random variable representing the output of the encoder for factor $i$, $i = 1, 2, \ldots, N$. Let $X_i^{(k)}$ be the $k$-th entry of $X_i$, which is a real-valued random variable. The normalized mutual information between $L_i$ and $X_j^{(k)}$ is defined as in [11]:

$$\hat{I}(L_i; X_j^{(k)}) = \frac{I(L_i; X_j^{(k)})}{H(L_i)}. \quad (24)$$

Then, the multi-dimensional mutual information between $L_i$ and $X_j$ is defined by taking the maximum of $\hat{I}(L_i; X_j^{(k)})$ over $k$:

$$\hat{I}(L_i; X_j) = \max_k \hat{I}(L_i; X_j^{(k)}). \quad (25)$$

One might argue that it is mathematically more meaningful to take the normalized mutual information between $L_i$ and the whole $X_j$ instead. But we found that, as the dimensionality of the code space increases, the number of samples required to accurately estimate $H(X_j)$ and $H(X_j|L_i)$ increases exponentially. And when the number of samples is insufficient, even a randomly initialized encoder would be incorrectly computed to have normalized mutual information close to one, which makes the evaluation meaningless. So we have to settle with our current definition.
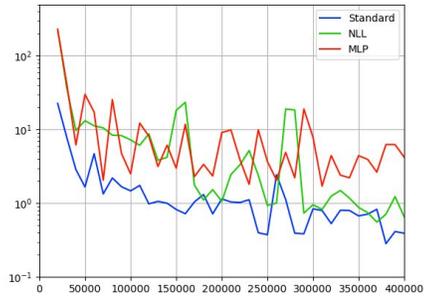


Figure 11: Average squared code distance.

Then, the MIG is computed as

$$\text{MIG} = \frac{1}{N} \sum_i (\hat{I}(L_i; X_i) - \max_{j \neq i} \hat{I}(L_i; X_j)). \quad (26)$$

As we have noted, in Table 2 in the main text we are only concerned with the disentanglement between the combined unknown factor and each individual labeled factor. To reflect this, in the computation of MIG here, in Equation 26 we only take the average over $i$ where factor $i$ is labeled.

Special procedures were taken to compute the MIG for [19]: the definition of MIG requires the distribution of $X_j$ to have continuous support, for otherwise all the normalized mutual information would be equal to one and the MIG would be zero. The output of [19] is in the form of an exact code, rather than a normal distribution as in VAE-based methods, so the $X_j$'s will have discrete support, making MIG non-applicable. Note that during the training of [19], Gaussian noise with fixed standard deviation was added to the content embedding, which effectively turns discrete codes into a distribution with full support. So in the computation of MIG, we similarly add Gaussian noise with a fixed standard deviation. The standard deviation is chosen using the following procedure: for all possible pairs of $(\sigma_1, \sigma_2)$ where $\sigma_1, \sigma_2 \in \{10^{-5}, 2 \times 10^{-5}, 5 \times 10^{-5}, 10^{-4}, \ldots, 1, 2, 5, 10\}$, we compute MIG by adding $\mathcal{N}(0, \sigma_1^2 I)$ to the content ("unknown factor" in our terminology) code and $\mathcal{N}(0, \sigma_2^2 I)$ to the class ("labeled factor") code. The values of $\sigma_1$ $\sigma_2$ are chosen so that the average MIG under two settings on the *3D chairs* dataset (rotation unknown and model unknown) is maximized. By this we determine that $\sigma_1 = 0.05$ and $\sigma_2 = 0.02$.

## C. Ablation Analysis

### C.1. Stability of Stage I

Our proposed methods of Negative Log Unlikelihood and sample-space classification aim to improve the stability of encoder-classifier adversarial training. Here we evaluate the effectiveness of these two schemes. Specifically,

(a) $\beta$-VAE [22]　　　　(b) Factor-VAE [29]　　　　(c) $\beta$-TCVAE [11]

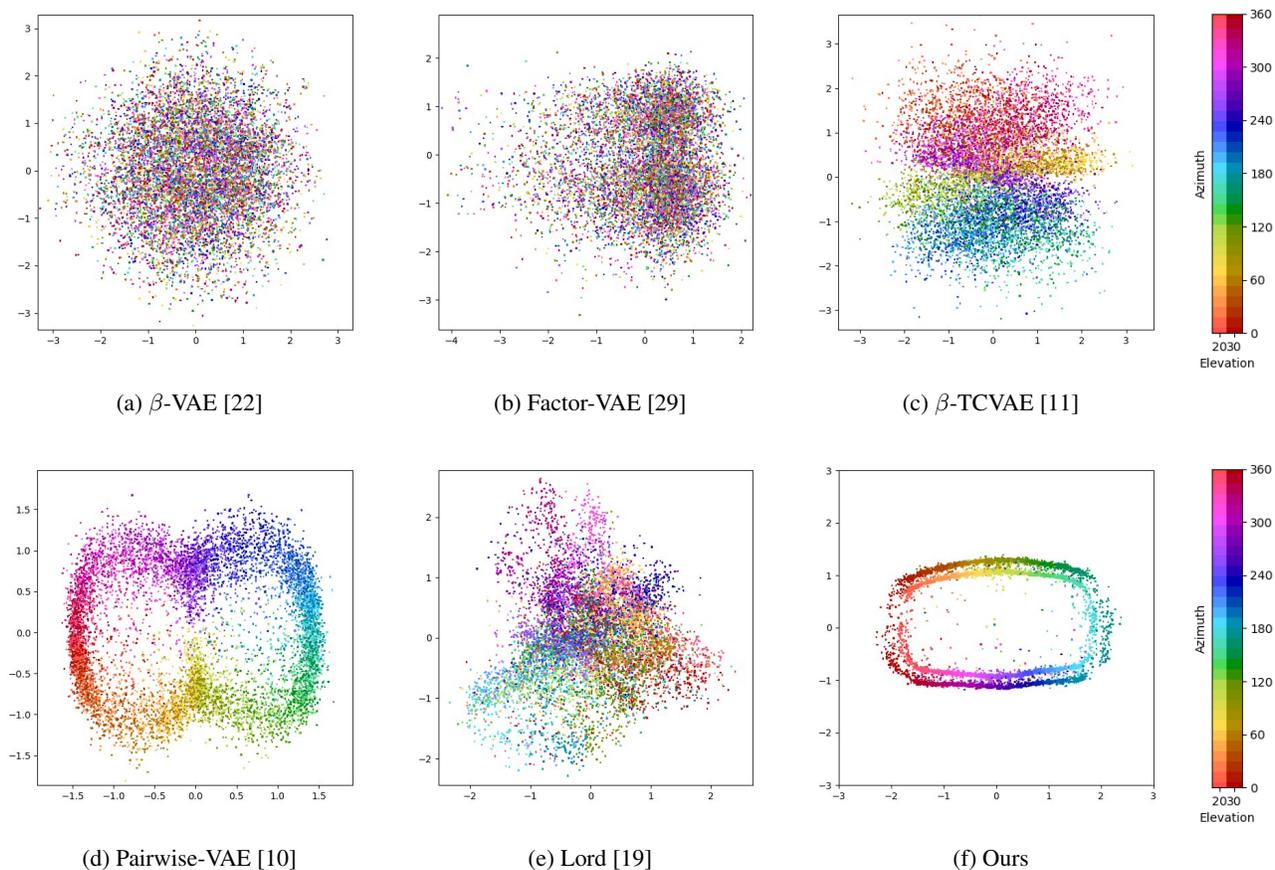(d) Pairwise-VAE [10]　　　　(e) Lord [19]　　　　(f) Ours

Figure 12: Visualizing the disentanglement with test sample distributions.

Table 6: Comparison of result with different Stage II configurations on MNIST.

| Configuration | MSE ↓ | MIG ↑ |
|---|---|---|
| Standard | **0.0086** | **0.978** |
| Non-adversarial $R$ | 0.0103 | 0.916 |
| No unknown code condition | 0.0402 | 0.930 |

we track the change of code distribution. We take snapshots of the network at fixed intervals during training. The whole test dataset is encoded, and we compute the average squared distance in the code space, from the code of each sample to the code of the same sample in the previous snapshot. In stable training, the encoder should keep the distance small while still finding a good distribution.

We train three variants of Stage I on the 3D chairs dataset with rotation unknown: one standard variant (proposed), one where the adversarial objective is maximizing the NLL loss of the classifier, and one where the classifier is an MLP operating in the code space, with four hidden layers of size
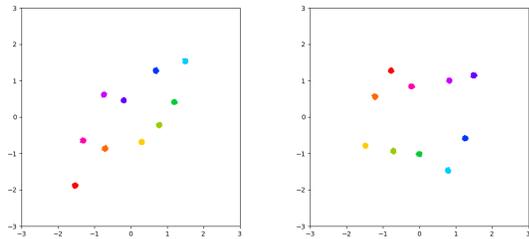
512.

The code distribution is computed every 10,000 iterations until iteration 400,000. The average squared code distance in the unknown code space every 10,000 iterations apart is shown in Figure 11, in logarithm scale. We can see that both NLU and sample-space classification contributed to reducing the fluctuation in code space.

## C.2. Adversarial Classifier and Condition on Unlabeled Code in Stage II

We evaluate the effectiveness of adversarial classifiers in Stage II compared to non-adversarial ones and examine the necessity of the code distance loss term. We train three variants of Stage II on the *MNIST* dataset: one standard variant (proposed), one without NLU term (remove NLU term from Equation 5e in the main text) so that the classifiers do not try to distinguish generated samples from real ones in the same class, and one without code distance term (remove code distance from Equation 5g in the main text) so that the network does not explicitly preserve the unknown factor.

The initial network weights of $E$ and $G$ for three configurations are inherited from the same Stage I training run

(a) Floor hue                    (b) Wall hue

Figure 13: Distribution of test samples in the code space of the two correlated labeled factors.



Figure 14: Distribution of test samples in the code space of the unknown encoder when the unknown factor and the labeled factors are correlated.

so that the result is only affected by Stage II training. We compute the mean squared reconstruction loss and mutual information gap for the final models in Table 6. By using the adversarial Stage II classifier and adding the code distance term, we are able to improve both disentanglement (MIG) and reconstruction (MSE).

## D. Effect of Factor Correlation

In the datasets used for evaluation, the factors are independent of each other. In particular, in the 3D Shapes and the 3D Chairs datasets, every combination of labels occurred exactly once. In this section we would like to explore the behavior of our method when some of the factors are correlated. For better control, we construct a dataset where the correlation is exactly known: in the 3D Shapes dataset, each of the three color factors has 10 possible values, numbered 0 to 9. We take the subset of the dataset consisting of all images whose *floor hue* and *wall hue* differ by 0 or $\pm 1$, modulo 10.

### D.1. Correlation Between Two Labeled Factors

The first case is when correlation exists between two of the labeled factors. The semantics of the labeled factors in our networks is enforced to strictly follow the semantics of the labels, so it is expected that our network would behave in the same way as when the labeled factors did not have correlation.

We train our model on the correlated subset with *object hue* being the unknown factor, so that the two correlated factors are both labeled. The distribution of test samples in the two correlated factors is shown in Figure 13. As can be seen, our network successfully learns to encode *floor hue* and *wall hue* as labeled.

While the behavior of the network remains the same, the MIG decreases: due to correlation the mutual information between *floor hue* and *wall hue* is now $\ln 10 - \ln 3$ instead of 0, and a perfect set of encoders would produce an MIG
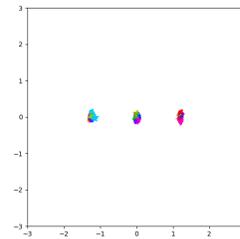
of

$$\frac{1}{6}\left(2\left(1 - \frac{\ln 10 - \ln 3}{\ln 10}\right) + 4\right) \approx 0.8257 \qquad (27)$$

In comparison, our method gives an MIG of $0.8026$.

### D.2. Correlation Between Labeled and Unknown Factors

The situation is more complicated when there is correlation between the unlabeled factor and the labeled factors. Our goal is for the unknown encoder to not encode any information about the labeled factors, which is to say, the conditional distribution of the unknown encoder, given the labeled factors, should be the same regardless of the value of the labeled factors. We train our model with *floor hue* being the unknown factor. In this case, since for any value of *wall hue* there are exactly three possible values of *floor hue*, it can be predicted that our unknown encoder should discover a factor with three discrete values, such that for any *wall hue*, each of the three *floor hues* is encoded as a different value. Note that this discovered factor is not necessarily the "hue difference" of the floor and the wall: there is no guarantee that the three values of this factor correspond to the hue difference being $-1$, $0$ and $1$ consistently, independent of other factors. The mapping from hue difference to the value of the discovered factor can vary according to *wall hue*. The only guarantee is that for the same *wall hue*, different *floor hues* correspond to different values of the discovered factor.

The distribution of the samples in the unknown encoder's code space, colored by *floor hue*, is shown in Figure 14. Three clusters can be clearly seen.

In general, we can conclude that if the intended semantics of the unknown factor is correlated to the labeled factors, then the factor discovered by our method would have different semantics. This may or may not be a desirable outcome, but it shows that our method is highly effective in ensuring the independence between the discovered unknown factor and the labeled factors.

(a) $\beta$-VAE [22]

(b) Factor-VAE [29]

(c) $\beta$-TCVAE [11]
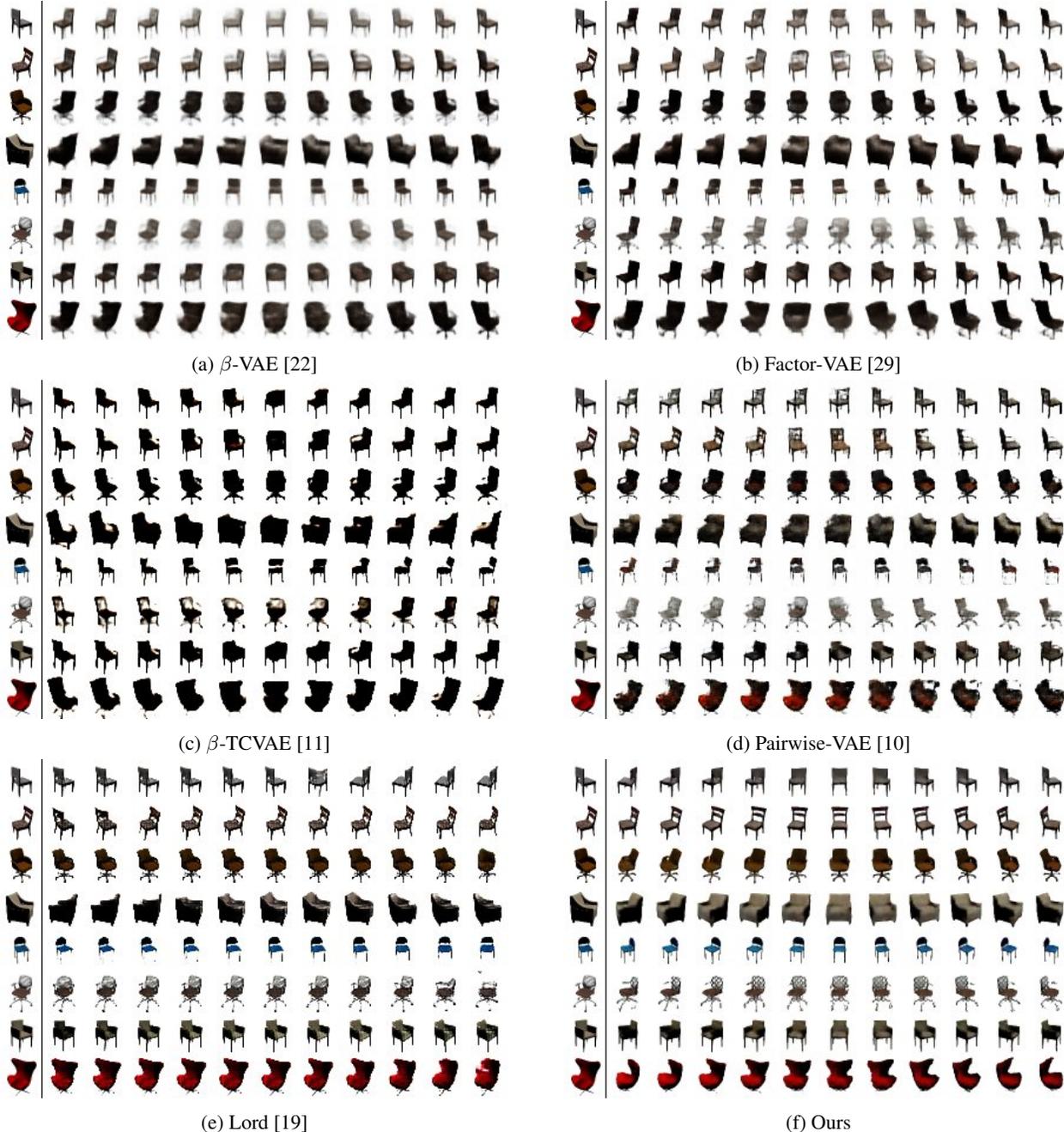
(d) Pairwise-VAE [10]

(e) Lord [19]

(f) Ours

Figure 15: Additional results of manipulation comparison on *3D Chairs* by uniformly sampling the latent codes depicting the azimuth rotation. The leftmost column shows the inputs.

# E. Comparisons

## E.1. Visualization

In Figure 12, for each of the methods compared, we plot the distribution of test samples of the *3D Chairs* dataset in the code space, projected onto the two dimensions having the largest mutual information with the rotation label. In $\beta$-VAE [22] and Factor-VAE [29], no clear color pattern can be observed. In $\beta$-TCVAE and Lord [19], samples with the same rotation are somewhat close together but there is no meaningful order between different rotations. Pairwise-VAE [10] and our method can arrange the azimuth angle correctly into a ring, but the structure is more clear in our method, and also, we can distinguish two slightly different elevation angles, which are not distinguished by any other method.

Figure 16: Additional portrait relighting results.

## E.2. Results

We provide more results of manipulating the latent code related to azimuth rotations and generate the images using different methods in Figure 15.

## F. Downstream Tasks

### F.1. Portrait Relighting

In this task, the labeled factor is lighting, represented by the coefficients of spherical harmonics up to second-order, which are 9 real-valued numbers. We show that our disentanglement framework can handle such continuous labels. More portrait relighting results are shown in Figure 16.

### F.2. Anime Style Transfer

Figure 17 and Figure 18 present more results of anime style transfer generated by fixing either style or content. In Figure 17, we show the styles of the results are consistently faithful to the input. In Figure 18, we explore the diversity of styles our network can learn and demonstrate the ability to generate the same content in different styles where facial shapes, appearances, and aspect ratios are captured.

### F.3. Landmark-Based Face Reenactment

The landmark-based face reenactment can generate the face motion of one target person from a single image with another source subject's facial expressions and head pose from a driving video. The generated results should preserve the source poses/expressions while matching the target identity. With only the identity label known, we can disentangle unknown poses/expressions from identities and synthesize new landmarks by changing identities in the source landmarks to the target.

We train the landmark-to-image network on the training dataset from VoxCeleb2 [15] which is processed by cropping a $256 \times 256$ face image and extracting its landmarks from each frame. In total, the training data contains $52,112$

Figure 17: Generated samples with fixed style and random content. Left column shows four input examples of two different artists.



Figure 18: Generated samples with fixed content and random style. The content code is fixed for all results while style codes randomly sample from the style distribution.

| Method | ID ↑ | Pose ↑ | Exp. ↓ |
|---|---|---|---|
| X2face [58] | 0.635 | 0.302 | 0.448 |
| First-order [52] | 0.770 | 0.822 | 0.274 |
| Ours *w/o* disentanglement | 0.715 | **0.862** | **0.208** |
| Ours *w/* disentanglement | **0.776** | 0.840 | 0.243 |

Table 7: Quantitative comparison of methods for cross-subject face reenactment on the VoxCeleb2 testing dataset between our method with and without landmark disentanglement, and two one-shot methods [58], and [52]. Note that our results without disentanglement are generated using the ground truth landmarks from the driving frame so their poses/expressions in the results should be the closest to those in the driving frame, but their identities are very different from the target person. Our method with disentanglement achieves the best in identity preservation and the second-best in poses/expressions.

videos for $1,000$ randomly selected subjects. We test on a video dataset of $8,000$ frames for 80 pairwise subjects (100 frames per video) randomly sampled from the VoxCeleb2 testing data,

We present more qualitative comparison in Figure 19 between with and without landmark disentanglement and in Figure 20 against two of the state-of-the-art one-shot face reenactment methods, *i.e.* X2face [58] and First-order [52]. Figure 19 compares the results using source landmarks without disentanglement and synthesized landmarks with disentanglement and shows that landmark disentanglement leads to better identity preservation. Figure 20 shows quantitative comparisons with two one-shot face reenactment baselines and demonstrates our method can better generalize to unseen subjects with more consistent quality under a large variety of poses/expressions.

Besides, we conduct a quantitative comparison as shown in Table 7 to measure the accuracy of disentanglement. We evaluate the results using three metrics (*ID*, *Pose*, and *Exp.*) for identity, pose, and expression respectively. *ID* measures the identity similarity between the resulting face and the target person by computing cosine similarity between their embedding vectors of the face recognition network VGGFace2 [7]. *Pose* measures the similarity between the resulting head pose and the source subject's pose in the driving frame by computing cosine similarity between their head rotations in radians around the X, Y, and Z axes estimated by OpenFace [3]. *Exp.* measures the difference between the resulting expression and the source expression in the driving frame by computing the L2 distance of their intensities of corresponding facial action units detected by OpenFace. Note that our results generated without disentanglement directly use the landmarks from the driving frame so their poses/expressions in the results are the closest

Table 8: Quantitative comparison with the state-of-the-art methods for skeleton disentanglement on Mixamo.

| Method | MSE ↓ | MAE ↓ |
|---|---|---|
| [63] | 0.0131 | 0.0673 |
| [1] | 0.0151 | 0.0749 |
| Ours | **0.0056** | **0.0498** |

to those in the driving frame but their identities mismatch the target person. In contrast, our method with disentanglement largely increases the identity accuracy while achieving comparable accuracy in poses/expressions.

### F.4. Skeleton-Based Body Motion Retargeting

We train our skeleton disentanglement network on the dataset of Mixamo [41]. It contains synthetic 3D skeletons of approximately $800$ unique motion sequences each of $36$ distinct characters ground truth labels for motion, view, and pose. We follow the same setting in [63] to use 32 of these characters with $800$ sequences of each one for training and the rest for testing. There is no overlap in motion sequence and character between training and testing. In training, we project the 3D joint coordinates on-the-fly onto 2D with viewing angles chosen randomly. The labeled factors are the identity of the character and the view angle and the unknown factor is the motion. The network sees only one frame at a time instead of a motion sequence.

In Table 8, we quantitatively compare with the state-of-the-art methods ( *i.e.*, [63] and [1]), which are task-specific, focusing on skeleton disentanglement. The method of [63] is unsupervised while the method of [1] utilizes full supervision. We perform evaluations on the same held-out test set from Mixamo (with ground truth available) using MSE and MAE as the metrics, reported in the original scale of the data. Our method, using weak supervision, outperforms both methods in terms of numerical joint position error but does not rely on any domain-specific prior knowledge.

Figure 21 shows that we can independently control identity, view and motion, synthesizing 2D skeletons with novel identity, view and motion while preserving the remaining factors unchanged.

Figure 22 shows more results of 2D skeleton-based motion retargeting. Although driven by different subjects, the target skeleton identity is completely preserved in the generated results.

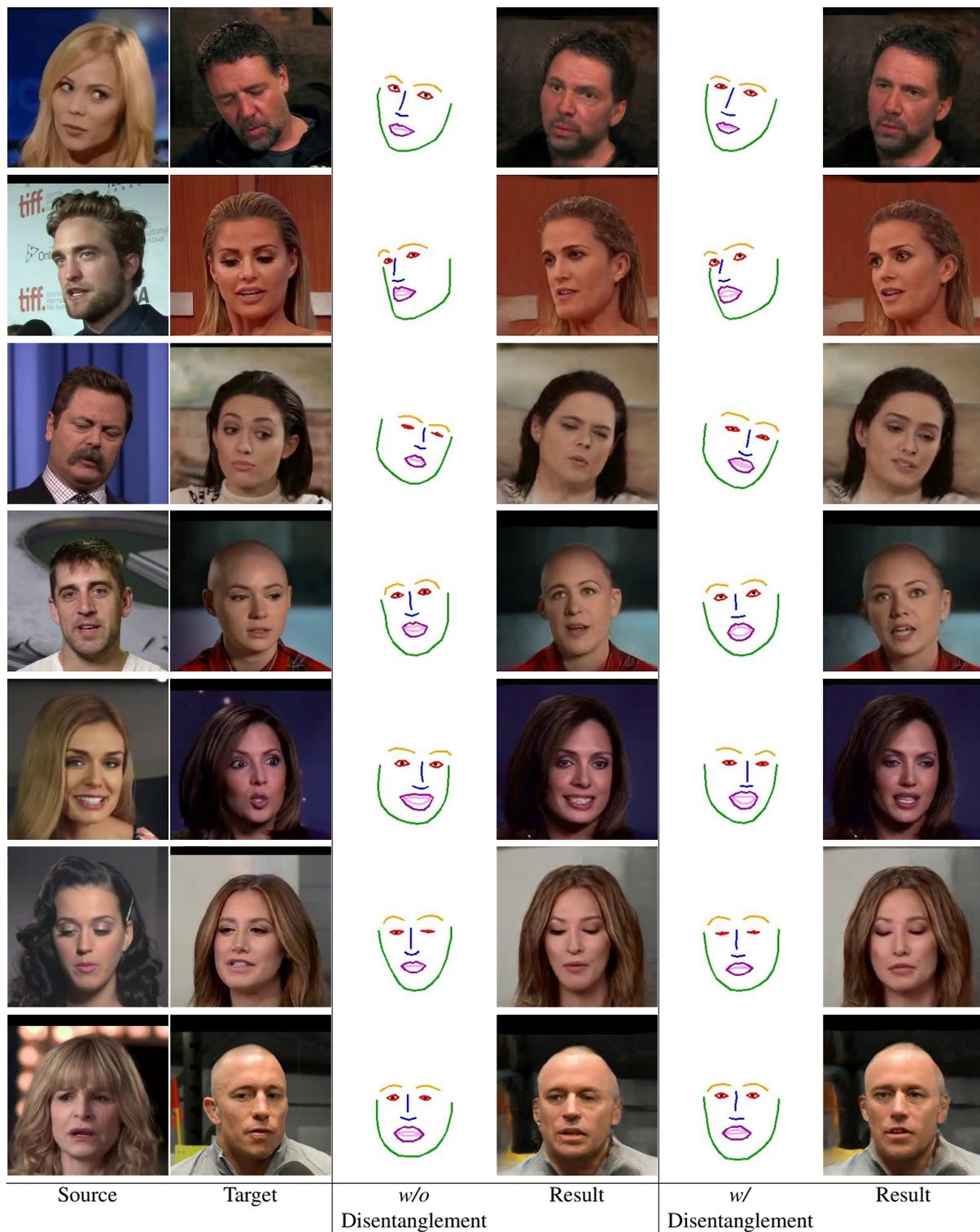| Source | Target | *w/o* Disentanglement | Result | *w/* Disentanglement | Result |
|--------|--------|--------|--------|--------|--------|

Figure 19: Qualitative comparison on face image reenactment between the translation results without (*w/o*) and with (*w/*) facial landmark disentanglement.

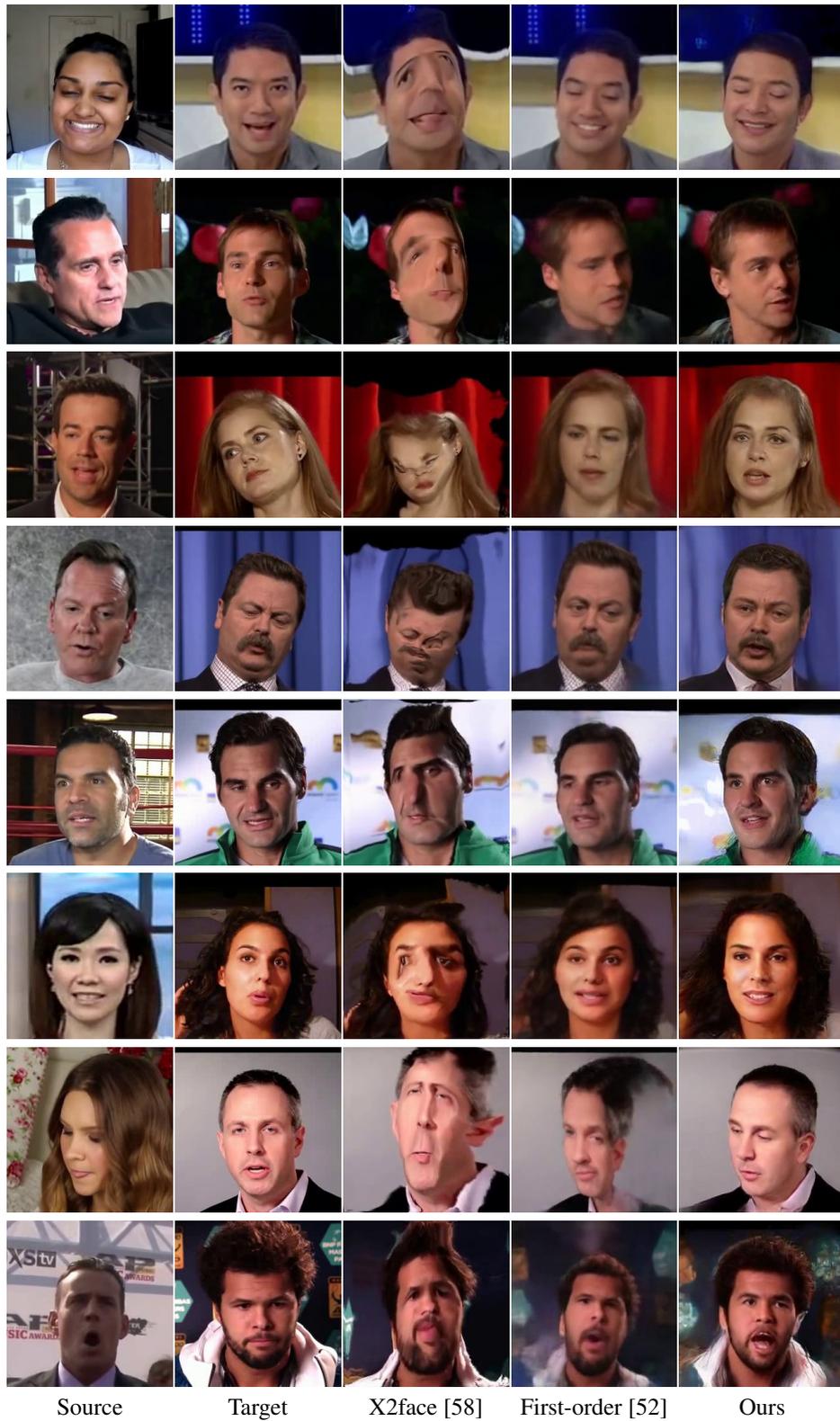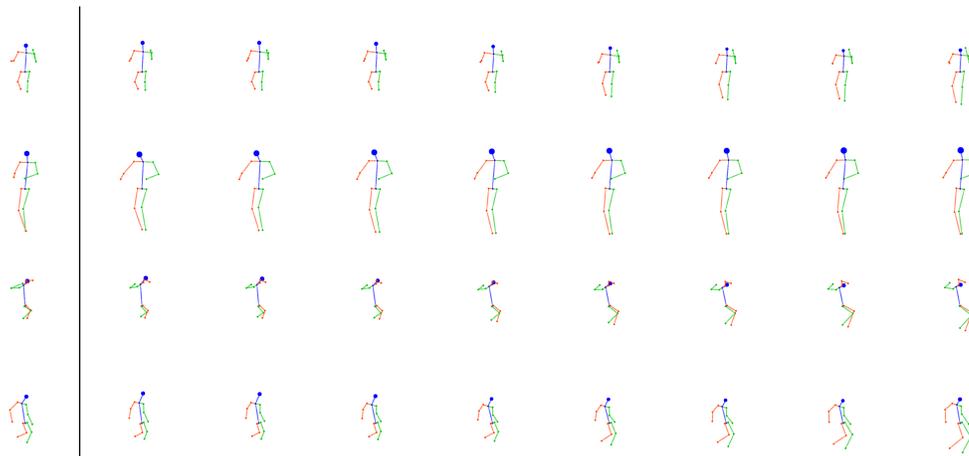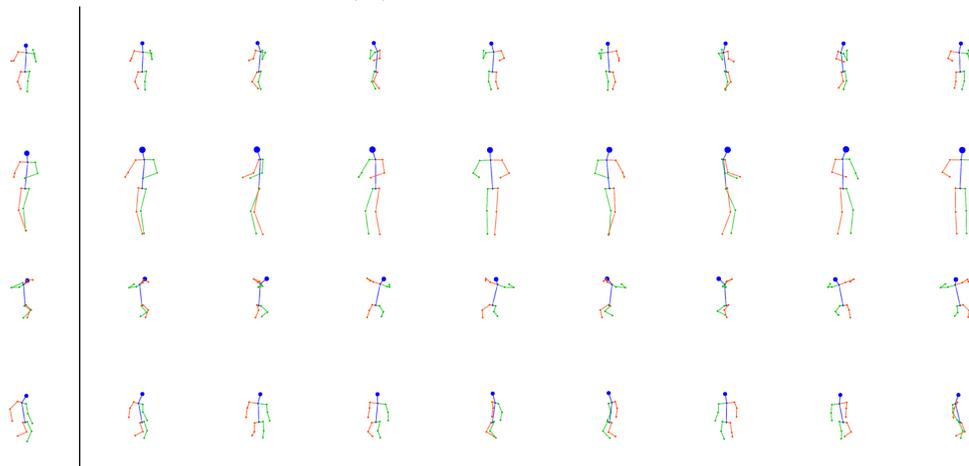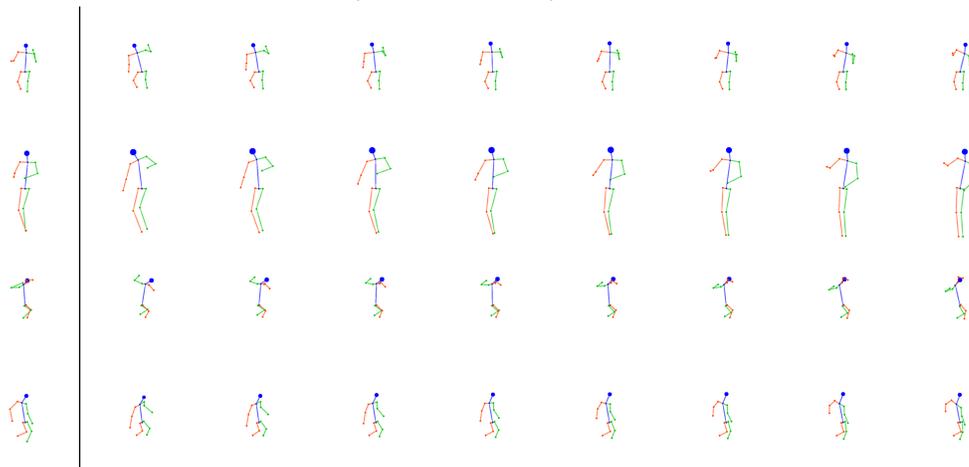| Source | Target | X2face [58] | First-order [52] | Ours |

Figure 20: Qualitative comparison on face image reenactment between our method and face motion transfer networks: X2face [58], and First-order [52] using images from Voxceleb2 [15] and FaceForensics++ [49].

(a) Identity synthesis with view and motion fixed


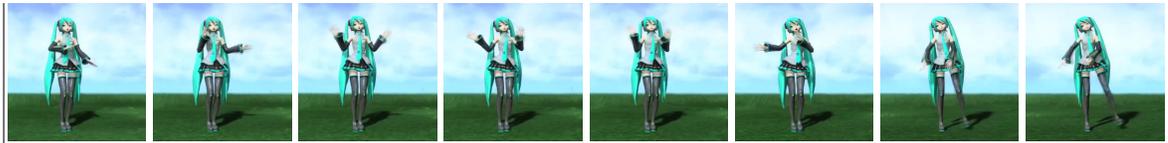
(b) View synthesis with identity and motion fixed



(c) Motion synthesis with identity and view fixed

Figure 21: Novel synthesis of identity, view and motion. Left column shows the input and the rest represents the generated skeletons by individually controlling identity, view and motion while fixing the others.

(a)



(b)

(c)



(d)

Figure 22: Motion retargeting results from different driving subjects to the same target person. In each case, left column shows the target person and his skeleton, right columns (from top to bottom) represent input source frames, extracted source skeletons, transformed skeletons, and generated target frames.