

# Escaping from Collapsing Modes in a Constrained Space

Anonymous ECCV submission

Paper ID 518

**Abstract.** Generative adversarial networks (GANs) often suffer from unpredictable mode-collapsing during training. We study the issue of mode collapse of Boundary Equilibrium Generative Adversarial Network (BEGAN), which is one of the state-of-the-art generative models. Despite its potential of generating high-quality images, we find that BEGAN tends to collapse at some modes after a period of training. We propose a new model, called *BEGAN with a Constrained Space* (BEGAN-CS), which includes a latent-space constraint in the loss function. We show that BEGAN-CS can significantly improve training stability and suppress mode collapse without either increasing the model complexity or degrading the image quality. Further, we visualize the distribution of latent vectors to elucidate the effect of latent-space constraint. The experimental results show that our method has additional advantages of being able to train on small datasets and to generate images similar to a given real image yet with variations of designated attributes on-the-fly.

## 1 Introduction

The main goal of this paper is to provide new insights into the problem of mode collapse in training Generative Adversarial Networks (GANs) [1]. GANs have shown great potential in generating new data based on real samples and have been applied to various vision tasks [2–9]. Our study points out a simple but effective approach that can be used to improve the stability of training GANs for generating high-quality images with respect to disentangled representations.

GANs comprise two core components: generator  $G$  and discriminator  $D$ . The two components are optimized with respect to two spaces. One is the latent space  $Z$  for the generator, and the other is the data space  $X$  associated with a real data distribution  $p_{\text{real}}(x)$  for training data  $x \in X$ . The objective of the generator is to find a mapping  $G : Z \rightarrow X$  that maximizes the probability of the discriminator mistakenly accepting a generated image  $G(z)$ ,  $z \in Z$  as from  $p_{\text{real}}(x)$ . On the contrary, the discriminator's objective is to distinguish whether any given  $x \in X$  belongs to  $p_{\text{real}}(x)$ . During training, the generator only learns from the information provided by the discriminator, and aims to estimate a good mapping such that  $p_{\text{model}}(G(z))$  is similar to  $p_{\text{real}}(x)$ .

Compared with auto-encoders [10], GANs can generate sharper images owing to the adversarial loss. However, a downside of adopting the adversarial loss is that it makes the training of GANs unstable. The performance is strongly

dependent on hyper-parameters selection, and the generated images tend to have weaker structural coherence.

Boundary Equilibrium Generative Adversarial Network (BEGAN) [11] introduced by Berthelot *et al.* suggests several modifications on the architecture and loss designs, which significantly improve the quality of generated images and the training stability. Another contribution of BEGAN is providing an approximation of convergence for the class of energy-based GANs.

Despite the promising improvements of BEGAN, we empirically observe that BEGAN still unavoidably runs into mode collapses after certain epochs of training. In our experiments, the exact time when mode collapsing happens is highly related to target image resolution and dataset size. In addition to the typical drawbacks of mode collapsing, this unpredictable behavior also makes BEGAN’s intended contribution to providing “global measure of convergence” incomplete.

## 1.1 Contributions

We propose a new constraint loss toward addressing the mode collapsing problem. We find that the mode-collapsing problem is suppressed after adding the constraint loss. This new loss term does not increase model complexity and is computationally low-cost. Furthermore, it does not introduce any trade-off regarding image quality and diversity. The proposed model is called *BEGAN with a Constrained Space* (BEGAN-CS).

We visualize the latent vectors produced in training phase using Principal Component Analysis (PCA) [12]. In section 3.1, we analyze the effect of the constraint loss and explain why this loss term makes training process stable.

Since BEGAN-CS is more stable during training, it performs consistently well even when the size of training dataset is ten-times smaller than the normal setting, in which BEGAN fails to obtain acceptable results. In section 4.3, our experiment shows that the proposed BEGAN-CS can eventually converge to a better state, while BEGAN ends up at mode collapsing in an early stage.

We further discover that BEGAN is able to learn strong and high-quality disentangled representations in an unsupervised setting. The learned disentangled representations could be used to modify the underlying attributes of generated images. In the meanwhile, owing to the constraint loss, BEGAN-CS can accomplish approximation  $Enc(x^*) \simeq z^*$  on-the-fly for any given real image  $x^*$ , where  $G(z^*)$  is an approximate image to  $x^*$  under the fixed generator weights. By leveraging the  $z^*$  approximation and the disentangled representations, BEGAN-CS can generate on the fly a set of images conditioning on a real image  $x^*$ . The generated images are visually similar to the given real image and are able to exhibit the adjustable disentangled attributes.

## 2 Related Work

Deep Convolutional Generative Adversarial Network (DCGAN) [13] improves the original GAN [1] by employing a convolutional architecture to achieve better

stability of training and enhanced quality of generated images. Salimans *et al.* further present several practical techniques for training GANs [14]. Nevertheless, avoiding mode collapsing while keeping the quality of generated images is still a challenging issue in practice.

Energy-Based Generative Adversarial Network (EBGAN) [15] introduces another perspective for formulating GANs. EBGAN implements the discriminator as an auto-encoder with per-pixel error. Boundary Equilibrium Generative Adversarial Network (BEGAN) [11] shares the same discriminator setting as EBGAN and makes several improvements on the designs of architecture and loss function. One of BEGAN’s core contributions is introducing the equilibrium concept, which balances the power between the generator and the discriminator. With these improvements, BEGAN provides fast and stable training convergence, and is capable of generating high visual-quality images. Another contribution of BEGAN is providing an approximate measure of convergence. The earlier class of GANs lacks convergence measurement. Not until later, a new class of GANs exemplified by Wasserstein Generative Adversarial Network (WGAN) [16] introduces a new loss metric, which correlates with the generator’s convergence. To our knowledge, BEGAN yields an alternative class of GANs that also has a loss correlated with convergence measurement.

Apart from the class of energy-based GANs, Progressive Growing of Generative Adversarial Networks (PGGANs) [17] is another approach to generating high-quality images. By changing the training procedure without modifying the original GAN loss, PGGANs are able to increase training stability and to produce diverse yet high-resolution (up to  $1024 \times 1024$  pixels) images.

### 3 Methods

Mode collapse is a phenomenon that the generated images get stuck in or oscillate between a few modes. This phenomenon under BEGAN’s setting has a unique characteristic. Since every sample shares the same encoder in the discriminator of BEGAN, the generated images that collapse at the same mode will share similar latent vectors as encoded by the encoder.

By leveraging this property, we propose the *latent-space constraint loss* ( $\mathcal{L}_c$ ), or the *constraint loss* for short. It constrains the L1-norm of the difference between the latent vector  $z$  and the internal state of encoder  $Enc(G(z))$ , where  $Enc$  is the encoder within the discriminator. During the training process, the constraint loss is only optimized with respect to the discriminator. Although the mode-collapsing problem happens on the generator side, adding the constraint loss directly to the generator would expose too much information to the generator about how to exploit the discriminator, and thus turns out accelerating the occurrence of mode collapse. The constraint loss can also be viewed as a regularizer, which guides the function  $Enc(G(\cdot))$  to be an identity function, and forces the encoder of the discriminator to retain the diversity and uniformity of randomly sampled  $z \in Z$ .

The objective function of BEGAN-CS is mostly similar to BEGAN, except the additional constraint loss. The full objective of BEGAN-CS includes

$$\mathcal{L}_G = \mathcal{L}(G(z_G; \theta_G); \theta_D), \quad \text{for } \theta_G \quad (1)$$

and

$$\mathcal{L}_D = \mathcal{L}(x_{\text{real}}; \theta_D) - k_t \cdot \mathcal{L}(G(z_D; \theta_G); \theta_D) + \alpha \cdot \mathcal{L}_c, \quad \text{for } \theta_D \quad (2)$$

with

$$\begin{cases} \mathcal{L}_c = \|z_D - Enc(G(z_D))\|_1, & \text{(the constraint loss)} \\ k_{t+1} = k_t + \lambda(\gamma \mathcal{L}(x; \theta_D) - \mathcal{L}(G(z_G; \theta_G); \theta_D)), & \text{for each epoch.} \end{cases} \quad (3)$$

The total loss  $\mathcal{L}_G$  of the generator and the total loss  $\mathcal{L}_D$  of the discriminator are optimized to solve for the parameters  $\theta_G$  and  $\theta_D$ , respectively. The function  $\mathcal{L}(x; \theta_D) = \|x - D(x)\|_1$  associated with  $\theta_D$  computes the L1-norm of the difference between any given image  $x$  and its reconstructed image  $D(x)$  by the decoder of the discriminator. The latent vectors  $z_D$  and  $z_G$  are randomly sampled from  $Z$ . The variable  $k_t \in [0, 1]$  controls how much emphasis to put on  $\mathcal{L}(G(z_D; \theta_G); \theta_D)$ . The hyper-parameter  $\gamma \in [0, 1]$  balances between the real-image reconstruction loss  $\mathcal{L}(x; \theta_D)$  and the generated-image discrimination loss  $\mathcal{L}(G(z_G; \theta_G); \theta_D)$ . The hyper-parameter  $\alpha$  is a weighting factor for constraint loss. The constraint loss  $\mathcal{L}_c$  is to enforce  $Enc(G(\cdot))$  to be an identity function for  $z_D$ .

### 3.1 Latent Space Analysis

For further illustrating the effectiveness of our method and analyzing the root cause of mode collapsing, we visualize the latent space through time with and without the constraint loss. We take PCA as our choice of dimensionality reduction method, and project the latent vectors onto two-dimensional space. Another common choice of dimensionality reduction for visualization is t-Distributed Stochastic Neighbor Embedding (t-SNE) [18]. For the latent space, we are more interested in the density and distribution of the points rather than the relative nearness between points or clusters. As a result, PCA is more suitable for our analysis.

Fig. 1 shows a preliminary analysis of BEGAN and BEGAN-CS. We train both models on the CelebA dataset [19]. The 64-dimensional latent vectors of generated images ( $Enc(G(z))$ ) and real images ( $Enc(x)$ ) are projected onto two-dimensional space via PCA.

In this experiment, BEGAN gets into mode collapse at epoch 23. In addition to the obvious change in the shape of distribution after BEGAN mode-collapsing, our empirical analysis also shows two strong patterns. First, in comparison with BEGAN, the latent-vector distribution (in red) of images generated by BEGAN-CS can better fit the real images' latent-vector distribution (in blue). The latent vectors of GEGAN-CS scatter more uniformly across all epochs.

Second, for BEGAN without adding the constraint loss, both the variance of real images' latent vectors ( $\text{Var}(\text{real})$ ) and the variance of generated images'

latent vectors ( $\text{Var}(\text{gen})$ ) grow rapidly as the number of epochs increases. Our hypothesis is that the latent spaces of real images and generated images both expand too rapidly and non-uniformly. Since the number of training data is fixed, as the latent space of real images expands, the density of real images decreases. In the end, the generator of BEGAN reaches a low-density area in the latent space where there is only a few latent vectors of real images nearby. The generator of BEGAN then gets stuck in that area. In contrast, BEGAN-CS has the latent-space constraint as a regularizer, which restricts the latent spaces of real images and generated images expand inadvertently. In other words, the constraint loss limits the distribution of  $\text{Enc}(G(z))$  to be similar to uniform distribution.

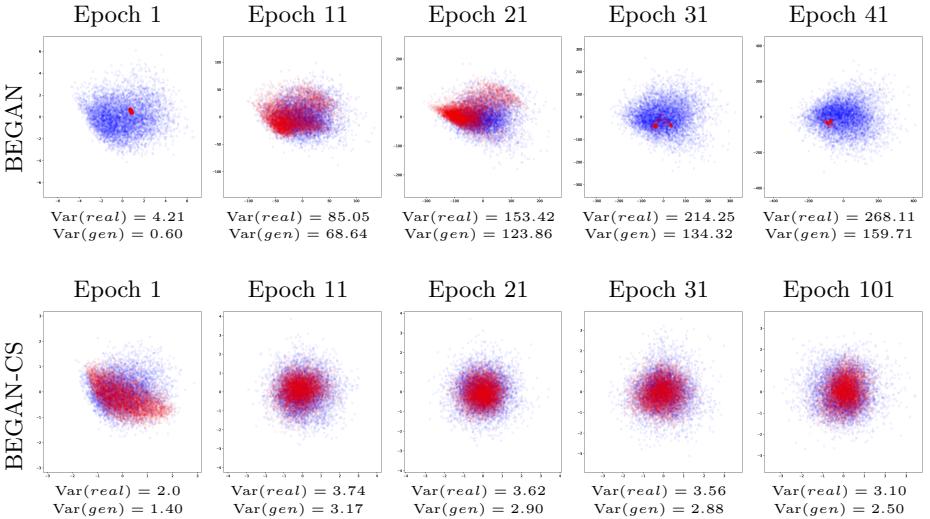


Fig. 1: We visualize the distributions of latent vectors of BEGAN and BEGAN-CS over epochs. Both BEGAN and BEGAN-CS are trained on CelebA dataset under  $64 \times 64$  resolution and batch size 64. Each graph consists of 6,400 random real images' latent vectors, *i.e.*  $\text{Enc}(x)$ , and 6,400 generated images' latent vectors, *i.e.*  $\text{Enc}(G(z))$ . The upper five graphs are generated by BEGAN, while the bottom five graphs are produced by BEGAN-CS. PCA is performed separately at each epoch based on the latent vectors of the real images. Each blue point represents a latent vector of a real image after applying PCA, and the red points correspond to the latent vectors of the generated images. The text under each graph lists the variance of real images' latent vectors ( $\text{Var}(\text{real})$ ) and the variance of generated images' latent vectors ( $\text{Var}(\text{gen})$ ). During the training of BEGAN, the variances of the distributions of latent vectors keep growing. Note that most of the graphs are created with a fixed interval of 10 epochs, except the bottom-right graph directly skips to the 101st epoch to highlight the effectiveness of BEGAN-CS. BEGAN has already collapsed before the 41st epoch.

### 225 3.2 Obtaining Optimal $z^*$ in One-Shot

226  
 227 Given an image  $x^*$ , finding an optimal latent vector  $z^*$  such that  $\|G(z^*) - x^*\|_1 <$   
 228  $\epsilon$  for some small  $\epsilon$  is a challenging problem for GANs. Traditionally,  $z^*$  can  
 229 be obtained by back-propagation for solving the optimization  $\min_{z^*}(\|G(z^*) -$   
 230  $x^*\|)$ . We name this optimization process as  $z^*$ -search. However,  $z^*$ -search is  
 231 time-consuming and needs to run for each inference individually, and thus is  
 232 impractical for real-world applications.  
 233

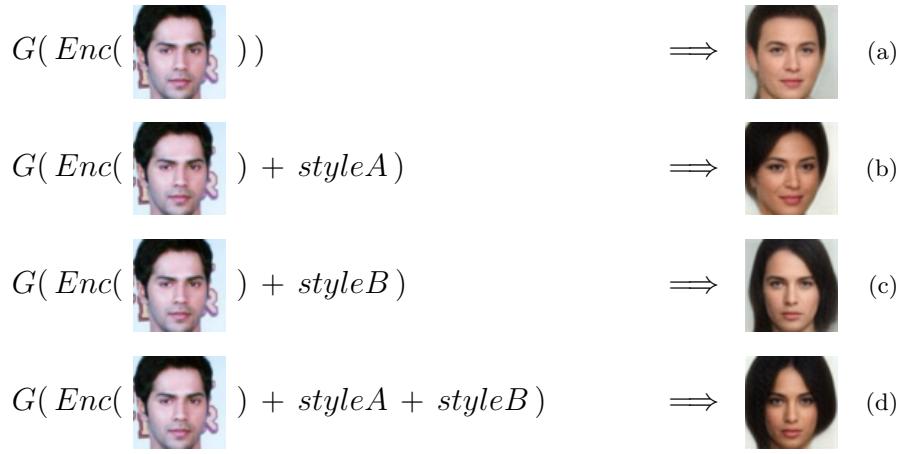
234 In the case of BEGAN-CS, the constraint loss works as a regularizer, guiding  
 235 the composite function  $Enc(G(z)) \simeq z$  to be similar to an identity function.  
 236 Consider the definition of  $z^*$ , where  $G(z^*) = x^*$ . We know that  $Enc(G(z^*))$   
 237 should be close to  $z^*$  due to the identity property. This implies that we may take  
 238  $x^*$  and obtain  $Enc(x^*)$  as an approximation to  $z^*$  after a single pass through  
 239 the encoder  $Enc(x^*)$ .  
 240

### 241 242 3.3 Disentangled Representation Learning and Application

243  
 244 We find that BEGAN is able to learn strong and high-quality disentangled rep-  
 245 resentations in an unsupervised setting. The direction of any vector within latent  
 246 space  $Z$  has a universally meaningful semantic, such as mixture of gender, age,  
 247 smile and hair-style. These learned representations can be combined with vector  
 248 arithmetic operations to generate images with multiple designated representa-  
 249 tions.  
 250

251 However, these disentangled representations are only effective for latent vec-  
 252 tors, which is a strong restriction that forbids many GAN models to use the  
 253 disentangled representation for practical applications, since obtaining the latent  
 254 vectors via  $z^*$ -search is computation-demanding. In the meanwhile, as we have  
 255 shown in section 3.2, BEGAN-CS is able to produce the approximation of  $z^*$  on-  
 256 the-fly. By adding multiple selected representation vectors to the approximated  
 257  $z^*$  with respect to any given real image  $x^*$ , we can generate images that are vi-  
 258 sually similar to  $x^*$  and comprise the selected representations at the same time.  
 259 We demonstrate this idea with a real example produced by BEGAN-CS in Fig 2.  
 260 In this example, the generated image of Fig 2d acquires both hair-styles shown  
 261 Fig 2b & Fig 2c. For BEGAN, which lacks the ability of estimating  $z^*$  directly,  
 262 the same effect may be forcibly achieved through time-consuming  $z^*$ -search to  
 263 obtain suitable  $z^*$ . Unfortunately,  $z^*$ -search causes the major bottleneck at in-  
 264 ference time and is therefore hard to use in real-world scenarios.  
 265

266 Similar applications can also be achieved using Variational Auto-Encoder  
 267 (VAE) based models [20–22] or other task-specific GAN models, such as In-  
 268 foGAN [23]. However, the images generated by VAE-based models tend to be  
 269 blurry, while InfoGAN cannot generate high-quality results as BEGAN does. In  
 comparison, our results are more promising in terms of stability and quality.



284 Fig. 2: An example of disentangled representations. The “*styleA*” and “*styleB*”  
 285 are two learned disentangled representations. Note that these representations are  
 286 universal and can be applied to any latent vector  $z$  for generating images  $G(z +$   
 287  $\text{style})$  with designated attributes. (a) Approximate  $z^*$  by  $G(\text{Enc}(x^*))$  in one-  
 288 shot. (b) & (c) The learned disentangled representations can be combined with  
 289  $G(\text{Enc}(x^*))$ . (d) Vector arithmetic with multiple disentangled representations. In  
 290 this case, the generated image has both hair-styles shown in *styleA* and *styleB*.  
 291  
 292

## 293 4 Experiments

296 We train BEGAN-CS using the CelebA dataset for all the experiments presented  
 297 in this paper. BEGAN-CS does not adopt the learning rate decay technique  
 298 described in BEGAN’s original paper, since the training process of BEGAN-CS  
 299 is already very stable. The hyper-parameter  $\alpha$  that controls the importance of  
 300 the constraint loss is set to 1 as the default value. For any hyper-parameter that  
 301 is not mentioned, we choose the same value as in BEGAN’s original setting.  
 302  
 303

### 304 4.1 Effectiveness of the Constraint Loss

306 In Fig. 3, we validate the effectiveness of the constraint loss. We show the gener-  
 307 ated images at specific epochs during the training of BEGAN and BEGAN-CS  
 308 on the CelebA dataset. The image resolution is  $64 \times 64$  and the batch size is  
 309 64. BEGAN-CS can continuously be trained up to 100 epochs without any evi-  
 310 dence of mode collapsing, loss of diversity, or reduction in quality. In contrast,  
 311 BEGAN encounters mode collapse at the 25th epoch (*i.e.*, the time-step B in  
 312 Fig. 3). In addition to the advantage of preventing from mode collapse, the pro-  
 313 posed BEGAN-CS model also maintains a very good performance in generating  
 314 high-quality images.

315  
316  
317  
318  
319  
320  
321  
322  
323  
324  
325

(A)



(B)



(C)

326  
327  
328  
329  
330  
331  
332  
333

BEGAN-CS



334

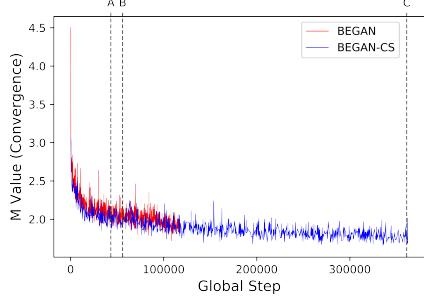
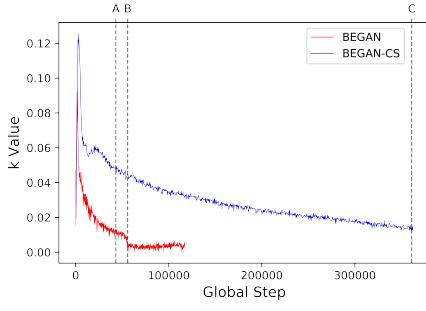
335  
336  
337  
338  
339  
340  
341  
342  
343  
344345  
346  
347  
348  
349  
350  
351  
352  
353  
354

Fig. 3: We validate the effectiveness of the constraint loss by showing the generated images at specific epochs during the training of BEGAN and BEGAN-CS on the CelebA dataset. The image resolution is  $64 \times 64$  and the batch size is 64. Note that BEGAN fails to reach epoch C since it already collapses at epoch B. In contrast, BEGAN-CS survives after epoch C. Furthermore, BEGAN-CS maintains a very good performance in generating high-quality images.

353  
354  
355  
356  
357  
358  
359  

## 4.2 Observing the Sudden Mode Collapsing

An interesting finding during our experiments is the timing of mode collapsing. As is mentioned in [11], the global measure of convergence can be used by BEGAN to determine whether the network has reached the final state or if the model has collapsed. However, in practice we are not able to observe significant evidence of mode collapsing directly from the value of the convergence measure.

315  
316  
317  
318  
319  
320  
321  
322  
323  
324  
325  
326  
327  
328  
329  
330  
331  
332  
333  
334  
335  
336  
337  
338  
339  
340  
341  
342  
343  
344  
345  
346  
347  
348  
349  
350  
351  
352  
353  
354  
355  
356  
357  
358  
359

360 Instead, the evidence of mode collapse are more often to be observed from the  $k$   
 361 value. The  $k$  value in BEGAN controls how much attention is paid on  $\mathcal{L}(G(z))$ .  
 362 According to our observation, every time the  $k$  value suddenly drops, BEGAN  
 363 is going to collapse shortly.

### 365 4.3 Better Convergence on Small Datasets

366  
 367 The dataset size is also an important factor for the timing of mode collapse.  
 368 Under a setting of reducing the training dataset CelebA to 1/10 of its original  
 369 size, BEGAN collapses earlier than training on full dataset. The early occurrence  
 370 of mode collapse keeps BEGAN from converging to an optimal state. The time-  
 371 step A in Fig. 4 is the best state that BEGAN can achieve during its training  
 372 on the down-sized CelebA dataset. On the other hand, BEGAN-CS has a more  
 373 stable training process. In Fig. 4, BEGAN-CS can continuously optimize on  
 374 the 1/10 down-sized CelebA dataset without encountering mode collapse, and  
 375 eventually converges to a better state than BEGAN.

### 376 4.4 Obtaining Optimal $z^*$ in One-Shot

377 In section 3.2, we have shown that BEGAN-CS can approximate optimal  $z^*$  with  
 378  $Enc(x^*)$ . Fig. 5 shows the experimental results of interpolation with obtained  $z^*$   
 379 from  $z^*$ -search using different GAN architectures. The experiments may serve as  
 380 proofs of concept for comparing the well-known GANs architectures, including  
 381 FisherGAN [24], PGGAN [17], and BEGAN. The experimental results show that  
 382 the obtained  $G(Enc(x^*))$  of BEGAN-CS is visually similar to  $x^*$ . In contrast,  
 383 the original BEGAN and other state-of-the-art GANs require time-consuming  
 384  $z^*$ -search for 10,000 iterations to obtain competitive results. It would take 340  
 385 seconds to 3,970 seconds depending on the network architecture. However, the  
 386 quality of the  $z^*$ -search result is still unstable and the searched image frequently  
 387 looks quite different to the given real image, such as wrong gender or incorrect  
 388 head pose. More examples on  $z^*$ -search with different GAN models and different  
 389 numbers of optimization iterations are shown in Fig. 6.

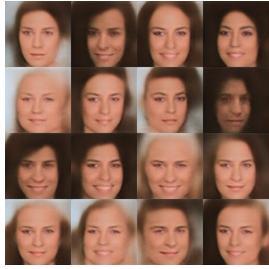
### 390 4.5 On-the-Fly Representation Manipulation

391 In section 3.3, we demonstrate a new application of BEGAN-CS with the dis-  
 392 entangled representations. By obtaining the approximation of  $z^*$  with  $Enc(x^*)$   
 393 and applying the selected disentangled representations, BEGAN-CS can gener-  
 394 ate images that are visually similar to  $x^*$  and exhibit the selected representations  
 395 at the same time. As a proof of concept, we visualize the process of adding single  
 396 representation in Fig. 7 and multiple representations in Fig. 9.

397 In Fig. 7, we first obtain the approximation of  $z^*$  from  $Enc(x^*)$ . Then for  
 398 each dimension  $i$ , we linearly interpolate and replace the the value of latent  
 399 vector  $z^*$  at its  $i$ th dimension by a grid value in  $[-5, 5]$  with step size 1, and  
 400 thus can generate a series of images based on the modified latent vectors. The

405  
406  
407  
408  
409  
410  
411  
412  
413  
414  
415  
BEGAN

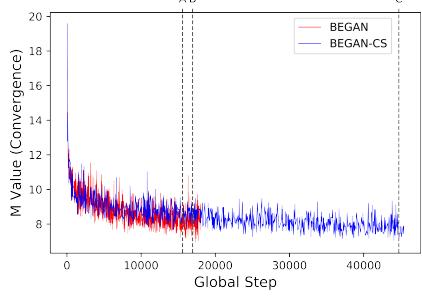
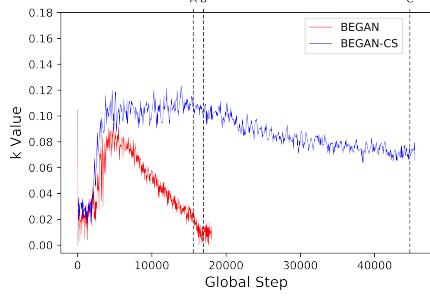
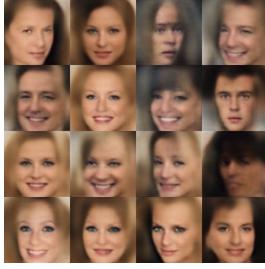
(A)



(B)



(C)

416  
417  
418  
419  
420  
421  
422  
423  
BEGAN-CS

435 Fig. 4: Better convergence of BEGAN-CS on small datasets. We show the generated images at selected epochs during training BEGAN and BEGAN-CS on 436 a 1/10 sized subset of CelebA. Training images are of  $128 \times 128$  resolution and 437 the batch size is 24. BEGAN-CS is stable and converges to a particularly better 438 state than BEGAN. The best state of BEGAN is at time-step A with degraded 439 quality, while BEGAN-CS can generate higher-quality results at time-step C. 440

441  
442  
443 images show that each dimension of the latent space  $Z$  represents a universal 444 disentangled representation. We can perform similar visual transformations to 445 any  $z \in Z$ . Fig. 7 shows some of the interesting disentangled representations. 446 The full visualization across the 64 dimensions is displayed in Fig. 8.

447 The learned disentangled representations can also be used to perform multi- 448 ple vector arithmetic operations on latent vectors. This property enables us 449 to control multiple attributes of a fixed image at the same time by adjusting

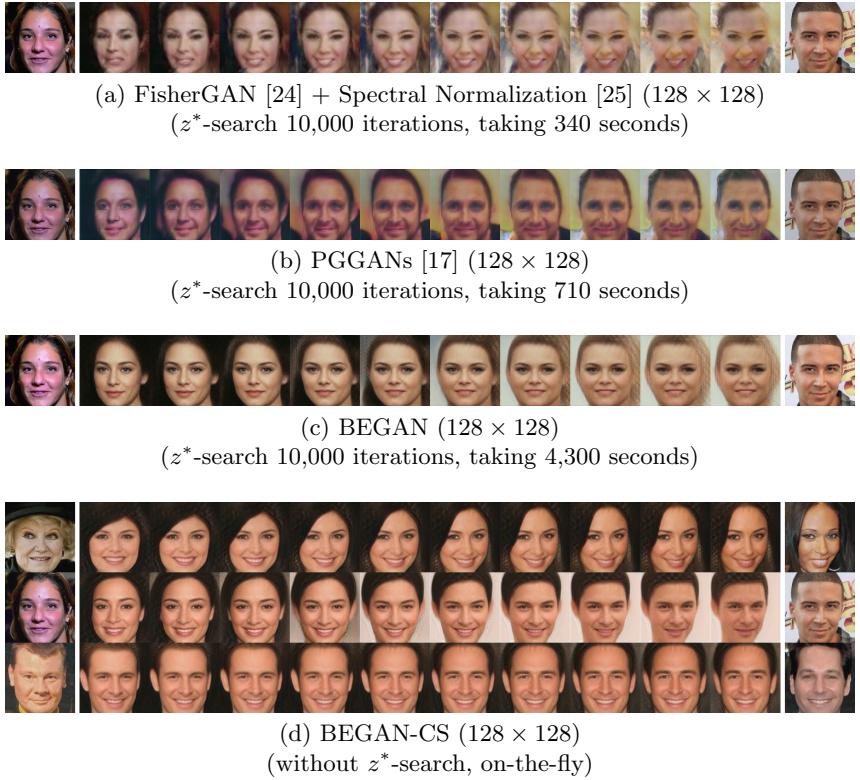


Fig. 5: Interpolation between two real images in latent space using different GAN models. Other state-of-the-art GANs require time-consuming  $z^*$ -search for 10,000 iterations to obtain competitive results, taking several minutes. Nevertheless, the quality of the  $z^*$ -search results is still not as good as the quality of the images generated on-the-fly by BEGAN-CS.

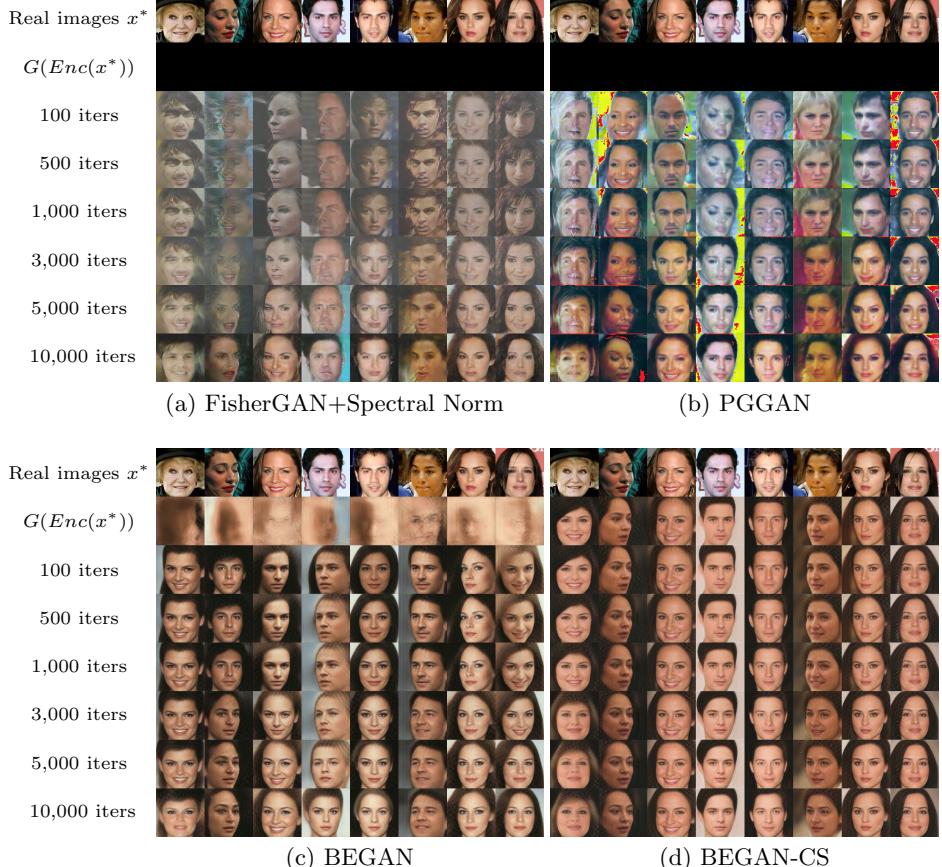
multiple dimension values on the corresponding latent vector. We visualize the results of combining two different representations in Fig. 9.

## 5 Conclusion

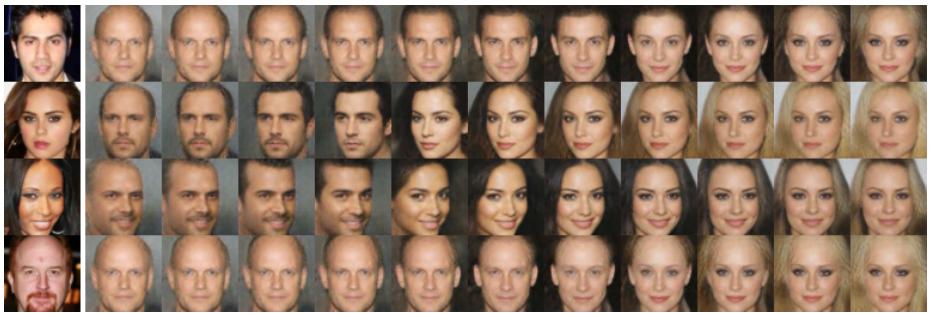
We identify that BEGAN suffers from the unpredictable mode-collapsing problem. The precise time when mode collapsing happens is non-deterministic, highly related to the resolution of generated images and the size of training dataset. We propose *BEGAN with a Constrained Space* (BEGAN-CS) toward addressing the mode-collapsing problem and visualize the effect of constraint loss in the latent space. We experimentally show that the model-collapsing problem is suppressed after adding the constraint loss. BEGAN-CS performs particularly better than BEGAN when the size of training dataset is ten-times smaller than the normal

495 setting. These advantages enable the class of energy-based GANs to move on to  
 496 the next challenge of generating even higher resolution images.

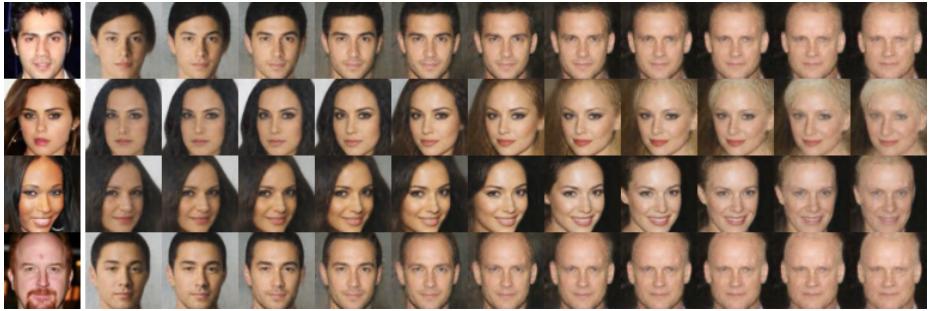
497 We also discover that BEGAN can learn salient and high-quality disentangled  
 498 representations in an unsupervised setting. Combined with the particular  
 499 property that BEGAN-CS is able to approximate  $z^*$  on-the-fly, BEGAN-CS can  
 500 generate images that are visually similar to the given real image and able to  
 501 exhibit the adjustable disentangled properties. “Obtaining  $z^*$  in one-shot” and  
 502 “adjustable image attributes” are two interesting properties that have various  
 503 potential applications, such as style manipulation and attribute-based editing.



535 Fig. 6: (a-c) We perform  $z^*$ -search from a random starting point  $z \in Z$  for  
 536 FisherGAN, PGGAN, and BEGAN. (d) BEGAN-CS starts from  $Enc(x)$ . We  
 537 also show the result of  $G(Enc(x^*))$  for BEGAN and BEGAN-CS. It can be seen  
 538 that only for BEGAN-CS the result of  $G(Enc(x^*))$  can be considered as a good  
 539 approximation of  $z^*$ .



(a) Gender.



(b) Age.



(c) Hair and skin color.

Fig. 7: Selected disentangled representations produced by BEGAN-CS at  $64 \times 64$  resolution. For each series of images, the left-most image is the fixed real image  $x^*$ . In each sub-figure, we first obtain approximation of  $z^*$  using  $Enc(x^*)$ . For each dimension  $i$ , we linearly interpolate and replace the  $i$ th dimension of  $z^*$  by a value in  $[-5, 5]$  with step size 1, and then generate the image set  $\{G(z_i^*)\}$ .

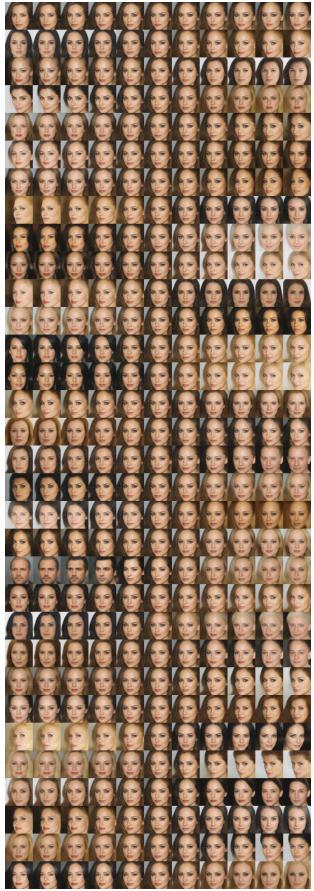
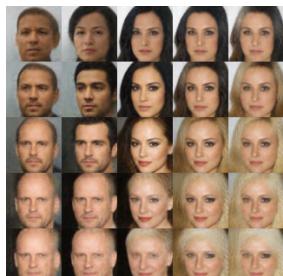


Fig. 8: Disentangled representations of BEGAN-CS across 64 dimensions along each axis in latent space  $Z$ .



(a) Row: gender. Column: hair and skin color.



(b) Row: gender. Column: age.



(c) Row: age. Column: hairstyle.

Fig. 9: Two-dimensional combinations of disentangled representations.

## 630 References

- 631 1. Goodfellow, I.J., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S.,  
Courville, A.C., Bengio, Y.: Generative adversarial nets. In: Advances in Neural  
Information Processing Systems 27: Annual Conference on Neural Information  
Processing Systems 2014, December 8-13 2014, Montreal, Quebec, Canada. (2014)  
2672–2680
- 632 2. Bousmalis, K., Silberman, N., Dohan, D., Erhan, D., Krishnan, D.: Unsupervised  
pixel-level domain adaptation with generative adversarial networks. In: 2017 IEEE  
Conference on Computer Vision and Pattern Recognition, CVPR 2017, Honolulu,  
HI, USA, July 21-26, 2017. (2017) 95–104
- 633 3. Dai, B., Fidler, S., Urtasun, R., Lin, D.: Towards diverse and natural image de-  
scriptions via a conditional GAN. In: IEEE International Conference on Computer  
Vision, ICCV 2017, Venice, Italy, October 22-29, 2017. (2017) 2989–2998
- 634 4. Gwak, J., Choy, C.B., Garg, A., Chandraker, M., Savarese, S.: Weakly supervised  
generative adversarial networks for 3d reconstruction. CoRR [abs/1705.10904](#)  
(2017)
- 635 5. Ledig, C., Theis, L., Huszar, F., Caballero, J., Cunningham, A., Acosta, A., Aitken,  
A.P., Tejani, A., Totz, J., Wang, Z., Shi, W.: Photo-realistic single image super-  
resolution using a generative adversarial network. In: 2017 IEEE Conference on  
Computer Vision and Pattern Recognition, CVPR 2017, Honolulu, HI, USA, July  
21-26, 2017. (2017) 105–114
- 636 6. Li, Y., Liu, S., Yang, J., Yang, M.: Generative face completion. In: 2017 IEEE  
Conference on Computer Vision and Pattern Recognition, CVPR 2017, Honolulu,  
HI, USA, July 21-26, 2017. (2017) 5892–5900
- 637 7. Shrivastava, A., Pfister, T., Tuzel, O., Susskind, J., Wang, W., Webb, R.: Learn-  
ing from simulated and unsupervised images through adversarial training. In: 2017  
IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2017, Honolulu,  
HI, USA, July 21-26, 2017. (2017) 2242–2251
- 638 8. Souly, N., Spampinato, C., Shah, M.: Semi supervised semantic segmentation using  
generative adversarial network. In: IEEE International Conference on Computer  
Vision, ICCV 2017, Venice, Italy, October 22-29, 2017. (2017) 5689–5697
- 639 9. Tzeng, E., Hoffman, J., Saenko, K., Darrell, T.: Adversarial discriminative do-  
main adaptation. In: 2017 IEEE Conference on Computer Vision and Pattern  
Recognition, CVPR 2017, Honolulu, HI, USA, July 21-26, 2017. (2017) 2962–2971
- 640 10. Hinton, G.E., Salakhutdinov, R.R.: Reducing the dimensionality of data with  
neural networks. science **313**(5786) (2006) 504–507
- 641 11. Berthelot, D., Schumm, T., Metz, L.: BEGAN: boundary equilibrium generative  
adversarial networks. CoRR [abs/1703.10717](#) (2017)
- 642 12. : Principal component analysis. Chemometrics and Intelligent Laboratory Systems  
**2**(1)
- 643 13. Radford, A., Metz, L., Chintala, S.: Unsupervised representation learning with  
deep convolutional generative adversarial networks. CoRR [abs/1511.06434](#)  
(2015)
- 644 14. Salimans, T., Goodfellow, I.J., Zaremba, W., Cheung, V., Radford, A., Chen, X.:  
Improved techniques for training gans. In: Advances in Neural Information Pro-  
cessing Systems 29: Annual Conference on Neural Information Processing Systems  
2016, December 5-10, 2016, Barcelona, Spain. (2016) 2226–2234
- 645 15. Zhao, J.J., Mathieu, M., LeCun, Y.: Energy-based generative adversarial network.  
CoRR [abs/1609.03126](#) (2016)

- 675 16. Arjovsky, M., Chintala, S., Bottou, L.: Wasserstein GAN. CoRR **abs/1701.07875** 675  
676 (2017) 676
- 677 17. Karras, T., Aila, T., Laine, S., Lehtinen, J.: Progressive growing of gans for im- 677  
678 proved quality, stability, and variation. CoRR **abs/1710.10196** (2017) 678
- 679 18. Maaten, L.v.d., Hinton, G.: Visualizing data using t-sne. Journal of machine 679  
680 learning research **9**(Nov) (2008) 2579–2605 680
- 681 19. Liu, Z., Luo, P., Wang, X., Tang, X.: Deep learning face attributes in the wild. In: 681  
682 Proceedings of International Conference on Computer Vision (ICCV). (2015) 681
- 683 20. Kingma, D.P., Welling, M.: Auto-encoding variational bayes. CoRR 683  
684 **abs/1312.6114** (2013) 684
- 685 21. Higgins, I., Matthey, L., Pal, A., Burgess, C., Glorot, X., Botvinick, M., Mohamed, 685  
686 S., Lerchner, A.: beta-vae: Learning basic visual concepts with a constrained vari- 685  
687 ational framework. (2016) 686
- 688 22. Larsen, A.B.L., Sønderby, S.K., Larochelle, H., Winther, O.: Autoencoding beyond 688  
689 pixels using a learned similarity metric. In: Proceedings of the 33nd International 689  
690 Conference on Machine Learning, ICML 2016, New York City, NY, USA, June 690  
691 19-24, 2016. (2016) 1558–1566 691
- 692 23. Chen, X., Chen, X., Duan, Y., Houthooft, R., Schulman, J., Sutskever, I., Abbeel, 692  
693 P.: Infogan: Interpretable representation learning by information maximizing gen- 693  
694 erative adversarial nets. In: Advances in Neural Information Processing Systems 694  
695 29: Annual Conference on Neural Information Processing Systems 2016, December 695  
696 5-10, 2016, Barcelona, Spain. (2016) 2172–2180 696
- 697 24. Mroueh, Y., Sercu, T.: Fisher GAN. In: Advances in Neural Information Processing 697  
698 Systems 30: Annual Conference on Neural Information Processing Systems 2017, 698  
699 4-9 December 2017, Long Beach, CA, USA. (2017) 2510–2520 699
- 700 25. Miyato, T., Kataoka, T., Koyama, M., Yoshida, Y.: Spectral normalization for 700  
701 generative adversarial networks. CoRR **abs/1802.05957** (2018) 701
- 702  
703  
704  
705  
706  
707  
708  
709  
710  
711  
712  
713  
714  
715  
716  
717  
718  
719