

MATCHINGGAN: MATCHING-BASED FEW-SHOT IMAGE GENERATION

Yan Hong, Li Niu, Jianfu Zhang, Liqing Zhang

ABSTRACT

To generate new images for a given category, most deep generative models require abundant training images from this category, which are often too expensive to acquire. To achieve the goal of generation based on only a few images, we propose matching-based Generative Adversarial Network (GAN) for few-shot generation, which includes a matching generator and a matching discriminator. Matching generator can match random vectors with a few conditional images from the same category and generate new images for this category based on the fused features. The matching discriminator extends conventional GAN discriminator by matching the feature of generated image with the fused feature of conditional images. Extensive experiments on three datasets demonstrate the effectiveness of our proposed method

Index Terms—Few-shot learning, generative adversarial network

1. INTRODUCTION

Deep generative models like Variational Auto-Encoder (VAE) [1] and Generative Adversarial Network (GAN) [2, 3] have made tremendous advance in generation problems. However, to generate new samples for a given category, the training of deep generative models relies on abundant labeled training data from this category. To achieve the goal of generation based on a few images, several few-shot generation methods [4, 5] have been proposed, which aim to generate more images from certain category based on a few images from this category. The model is trained on seen categories in the training stage, and then applied to generate more images for unseen categories which do not appear in the training stage. In this way, meta information or metric information could be transferred from seen categories to unseen categories, which makes few-shot generation possible for unseen categories. Existing few-shot generation methods can be categorized into feature generation methods [6] and image generation methods [7]. Specifically, given a few images (*resp.*, features) from one category, few-shot image (*resp.*, feature) generation methods target at generating more images (*resp.*, features) from this category.

In this paper, we focus on few-shot image generation, which is much more challenging than few-shot feature generation. Generated images can augment training data and facilitate downstream tasks like few-shot classification. To the

best of our knowledge, there are quite limited works [4, 5, 7] on few-shot image generation. For example, FIGR [4] incorporates meta-learning [8] into GAN to learn the data distribution of a few images from the same category. However, the quality of generated images is inferior and the model needs to be fine-tuned on the images from unseen categories at the test phase. GMN [5] combines matching procedure and VAE to generate novel images conditioned on a few images. However, due to the weakness of VAE, the images generated by GMN are vague and unrealistic. DAGAN [7] was designed to generate diverse images based on a single conditional image by injecting random noise into decoder. However, the diversity of generated images brought by injected noise is quite limited. Besides, DAGAN is conditioned on a single image, and thus fails to exploit the relationship among multiple images from the same category.

Considering the drawbacks of the above few-shot image generation methods, we propose Matching-based Generative Adversarial Network (MatchingGAN), inspired by matching-based methods [9, 5] which prove that matching procedure learned from training samples of seen categories can be adapted to unseen categories. Our MatchingGAN can fully exploit multiple conditional images from the same category to generate more diverse and realistic images by virtue of the combination of adversarial learning and matching procedure. In detail, our MatchingGAN is comprised of a matching generator and a matching discriminator. In the matching generator, we project a random vector and a few conditional images from one seen category into a common matching space, and calculate the similarity scores, which are used as interpolation coefficients to fuse features of conditional images to generate new images belonging to this seen category. Through the matching procedure, we can learn reasonable interpolation coefficients which determine how much information we should borrow from each conditional image. In the matching discriminator, we not only distinguish real images from fake images as done by conventional GAN discriminator, but also match the discriminative feature of generated image with the fused discriminative feature of conditional images to ensure that the generated image contains the interpolated information of conditional images. In the test phase, given a few conditional images from one unseen category, we can feed sampled random vectors into matching generator to generate numerous diverse and realistic images for this unseen category.

Our major contributions can be summarized as follows:

1) We propose a novel few-shot image generation method by combining matching procedure and adversarial learning; 2) Technically, we design a matching generator and a matching discriminator for our MatchingGAN; 3) Comprehensive generation and classification experiments on three datasets demonstrate the effectiveness of our MatchingGAN.

2. RELATED WORK

Generative adversarial network Generative Adversarial Network (GAN) [2] was proposed to discriminate real samples from fake samples and generate more realistic samples. In the early stage, unconditional GANs [3, 10] use random vectors to generate realistic samples based on the learned distribution of training samples. Then, GANs conditioned on a single image [11, 12, 7] were proposed to transform the conditional image to a target image by using adversarial learning. Recently, a few conditional GANs attempted to leverage more than one image to accomplish more challenging tasks, such as few-shot image translation [13] and few-shot image generation [4]. In this paper, we choose to combine matching procedure and GAN to capture the variance of a few conditional images to generate diverse images.

Few-shot learning Existing few-shot learning methods are mainly classified into three categories: metric-based methods [9, 14, 15], model-based methods [16], and optimization-based methods [17, 8, 18]. Our work is most related to metric-based methods [9, 14, 15], which can learn well-tailored metric to classify samples of unseen categories by comparing with a few labeled images from unseen categories. The process of learning metric well-tailored for specific tasks is defined as matching procedure.

Our MatchingGAN employs matching procedure to learn reasonable interpolation coefficients to fuse the features of conditional images.

Data augmentation Data augmentation [19] aims to generate more samples based on the given samples. Early data augmentation tricks such as shifts, rotations or shears can only produce limited diversity. In contrast, deep generative models could generate more diverse samples for data augmentation, including feature augmentation [20, 21, 6] and image augmentation [7]. Our method can be deemed as an image augmentation method. For image augmentation, previous few-shot image generation methods like FIGR[4], GMN [5], DAGAN [7] attempted to generate images based on a single or a few conditional images.

Besides, few-shot image translation method like FUNIT [13] intended to translate images from seen categories to unseen categories. However, in the testing phase, few-shot image translation relies on the images from seen categories to generate new images for unseen categories, which is a different task from few-shot generation. In this work, we propose a novel few-shot image generation method which could be used for data augmentation.

3. OUR MATCHINGGAN

3.1. Overview

Our MatchingGAN aims to learn a mapping from a few conditional images $\mathcal{X}_S = \{\mathbf{x}_i\}_{i=1}^K$ within one category to a new image $\tilde{\mathbf{x}}$ of this category, in which K is the number of conditional images. Let \mathcal{L}^s and \mathcal{L}^u be the set of seen categories and unseen categories respectively, where $\mathcal{L}^s \cap \mathcal{L}^u = \emptyset$. In the training phase, MatchingGAN is trained on images from seen categories \mathcal{L}^s , without reaching images from unseen categories \mathcal{L}^u . The trained model owns the ability to fuse the information of conditional images to generate new images from the same category.

In the testing phase, given K conditional images from an unseen category, the trained model could produce diverse and plausible images for this unseen category without any further fine-tuning. As illustrated in Fig. 1, our MatchingGAN consists of a matching generator and a matching discriminator, which will be detailed next.

3.2. Matching Generator

In our matching generator, there are three encoders including E_z , E_ϕ , and E_ψ as well as a decoder G . Encoders E_z and E_ϕ contribute to the matching procedure. Specifically, E_z (*resp.*, E_ϕ) projects a random vector \mathbf{z} sampled from unit Gaussian distribution $\mathcal{N}(\mathbf{0}, \mathbf{1})$ (*resp.*, conditional image \mathbf{x}_i) into a common matching space, leading to $E_z(\mathbf{z})$ (*resp.*, $E_\phi(\mathbf{x}_i)$). In the matching space, we calculate the similarity score between random vector \mathbf{z} and each conditional image \mathbf{x}_i as

$$a(E_z(\mathbf{z}), E_\phi(\mathbf{x}_i)) = \frac{e^{\cos(E_z(\mathbf{z}), E_\phi(\mathbf{x}_i))}}{\sum_{i=1}^K e^{\cos(E_z(\mathbf{z}), E_\phi(\mathbf{x}_i))}}, \quad (1)$$

in which $\cos(\cdot, \cdot)$ is the cosine similarity. We calculate the normalized cosine similarity as similarity score, which are later used as interpolation coefficients to fuse the features of conditional images. Through this matching procedure, we tend to learn reasonable interpolation coefficients to determine how much information we should borrow from each conditional image. Moreover, we expect that the matching procedure learned from seen categories can be adapted to unseen categories to generate images for unseen categories in the testing phase. Compared with simply using random interpolation coefficients without matching procedure, we verify that learning reasonable interpolation coefficients is more effective in our experiments (see Supplementary).

Encoder E_ψ and decoder G form an auto-encoder, which is used to generate new images based on the fused features of conditional images. Our encoder E_ψ and decoder G are designed as UNet structure [22], in which the features of several blocks in encoder E_ψ are connected to the output of corresponding blocks in the decoder G . Specifically, both E_ψ and G have four blocks, and we add skip connections between

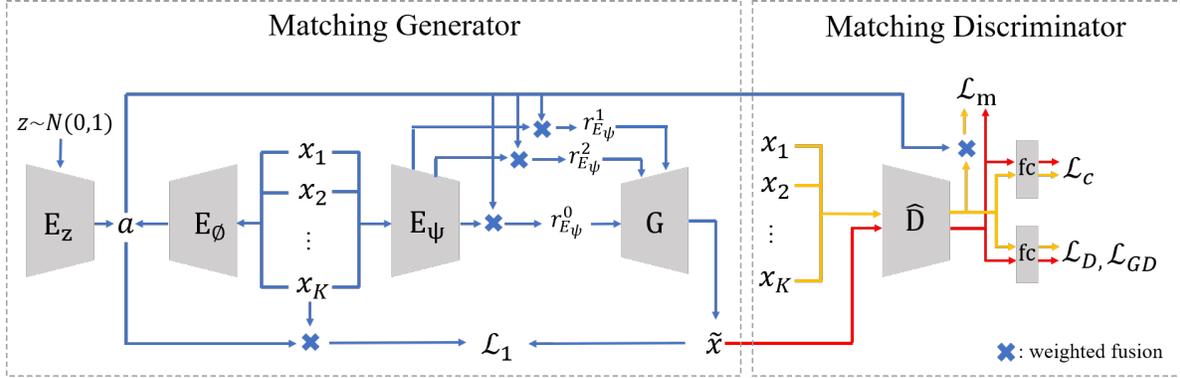


Fig. 1. The framework of our MatchingGAN which consists of a matching generator and a matching discriminator. \tilde{x} is generated based on the random vector z and K conditional images $\{x_i\}_{i=1}^K$. Best viewed in color.

the final two blocks in E_ψ and the first two blocks in G . We use $E_\psi^0(x_i)$ to denote the bottleneck feature of conditional image x_i , and use $E_\psi^j(x_i)$ to denote the feature of the j -th encoder block with skip connection. After using similarity scores $a(E_z(z), E_\phi(x_i))$ to interpolate the features $E_\psi^j(x_i)$ of K conditional images, we can have the fused features:

$$r_{E_\psi}^j = \sum_{i=1}^K a(E_z(z), E_\phi(x_i)) E_\psi^j(x_i), \quad j = 0, \dots, L, \quad (2)$$

in which L is the number of skip connections.

Due to the skip connection between encoder E_ψ and decoder G , we concatenate the fused feature $r_{E_\psi}^j$ from E_ψ^j and the output of its connected block in G as the input to the next block in G . Then, the image generated by the decoder G can be represented by

$$\tilde{x} = G(r_{E_\psi}^0, r_{E_\psi}^1, \dots, r_{E_\psi}^L). \quad (3)$$

As the shallow (*resp.*, deep) layer in the encoder integrates the low-level (*resp.*, high-level) information, our generated images can fuse multi-level information of conditional images coherently based on the interpolated features in multiple layers. The effectiveness of fusing multi-level features is proved in our experiments (see Supplementary).

Considering that the generated image \tilde{x} fuses the information of K conditional images x_i based on similarity scores $a(E_z(z), E_\phi(x_i))$, we conjecture that \tilde{x} should appear more similar to the conditional image with higher similarity score. Thus, we employ the weighted reconstruction loss as follows,

$$\mathcal{L}_1 = \sum_{i=1}^K a(E_z(z), E_\phi(x_i)) \|x_i - \tilde{x}\|_1. \quad (4)$$

The weighted reconstruction loss can make the network training more stable and enforce the generated images to contain the fused information of conditional images as expected.

3.3. Matching Discriminator

In our matching discriminator D , we treat K conditional images x_i as real images while the generated images \tilde{x} as fake images. In detail, we calculate the average of scores $D(x_i)$ for K conditional images x_i and score $D(\tilde{x})$ for the generated image \tilde{x} . To stabilize adversarial learning, we adopt the hinge adversarial loss in [11]. Concretely, the discriminator D tends to minimize the loss function \mathcal{L}_D while the matching generator tends to minimize the loss function \mathcal{L}_{GD} :

$$\begin{aligned} \mathcal{L}_D &= \mathbb{E}_{\tilde{x}}[\max(0, 1 + D(\tilde{x}))] + \mathbb{E}_{x_i}[\max(0, 1 - D(x_i))], \\ \mathcal{L}_{GD} &= -\mathbb{E}_{\tilde{x}}[D(\tilde{x})]. \end{aligned} \quad (5)$$

Following ACGAN [23], we construct classifier C by replacing the last fully connected (fc) layer of the discriminator D with another fc layer with C^s outputs, in which C^s is the number of seen categories. Then, we employ the cross-entropy classification loss to distinguish different categories:

$$\mathcal{L}_c = -\log p(c(x)|x), \quad (6)$$

where $c(x)$ is the category of image x . We train the discriminator by minimizing the classification loss $\mathcal{L}_c^D = -\log p(c(x_i)|x_i)$ of conditional images \mathcal{X}_S . When updating the matching generator, we expect the generated image \tilde{x} to be classified as the same category of conditional images by minimizing the classification loss $\mathcal{L}_c^G = -\log p(c(\tilde{x})|\tilde{x})$.

In practice, the distributions of real and fake images may not overlap with each other, especially at the early stage of training process. Hence, the discriminator D can separate them perfectly, which makes the training process of matching generator unstable. Considering the fact that the matching generator fuses the features of conditional images \mathcal{X}_S according to the interpolation coefficients, to cooperate with the fusion strategy in the matching generator, we match the discriminative feature of \tilde{x} with the fused discriminative feature of \mathcal{X}_S , which is implemented by a feature matching loss. In

detail, we remove the last fc layer of the discriminator D and use the remaining feature extractor \hat{D} to extract the discriminative features of generated images and conditional images. Thus, the feature matching loss can be written as

$$\mathcal{L}_m = \left\| \sum_{i=1}^K a(E_z(z), E_\phi(x_i)) \hat{D}(x_i) - \hat{D}(\tilde{x}) \right\|_1. \quad (7)$$

3.4. Optimization

The total loss function our MatchingGAN can be written as

$$\mathcal{L} = \mathcal{L}_D + \mathcal{L}_{GD} + \lambda_r \mathcal{L}_1 + \mathcal{L}_c + \lambda_m \mathcal{L}_m. \quad (8)$$

During adversarial learning, matching generator and matching discriminator are optimized by different loss terms in an alternating manner. In particular, the matching discriminator D is trained with \mathcal{L}_D and \mathcal{L}_c^D , while the matching generator E_z , E_ϕ , E_ψ , and G are trained with \mathcal{L}_{GD} , \mathcal{L}_1 , \mathcal{L}_c^G , and \mathcal{L}_m .

4. EXPERIMENTS

In this section, we compare our MatchingGAN with existing methods by conducting generation and classification experiments on three datasets.

4.1. Datasets and Implementation Details

Following DAGAN [7], we conduct experiments on three datasets: Omniglot [24], EMNIST [25], and VGGFace [26]. We also follow the category split used in [7]. For Omniglot (*resp.*, EMNIST, VGGFace), a total of 1623 (*resp.*, 48, 2395) categories are split into 1200 (*resp.*, 28, 1802) seen categories, 212 (*resp.*, 10, 497) validation seen categories, 211 (*resp.*, 10, 96) unseen categories. Validation seen categories are used to monitor the training procedure, but not engaged in updating model parameters. For EMNIST and VGGFace, some categories have more than 100 samples. Following [7], for these categories, we randomly choose 100 images from each category to fit a low-data setting.

We empirically set $\lambda_r = 0.1$ and $\lambda_m = 1$. The number of conditional images K is set as 3 considering the trade-off between effectiveness and efficiency. We use Adam optimizer with learning rate 0.0001 and train our MatchingGAN for 200 epochs. The detailed architecture of matching generator and matching discriminator is provided in Supplementary.

4.2. Quantitative Evaluation of Generated Images

We evaluate the quality of images generated by different methods on VGGFace dataset based on commonly used Inception Scores (IS) [29] and Frchet Inception Distance (FID) [30]. The implementation details of IS and FID are described in Supplementary.

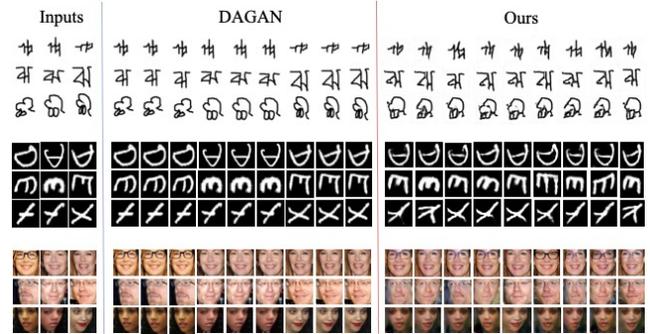


Fig. 2. Images generated by our MatchingGAN ($K = 3$) and DAGAN on three datasets (from top to bottom: Omniglot, EMNIST, VGGFace). The conditional images are in the left three columns.

Table 1. FID (\downarrow) and IS (\uparrow) of images generated by different methods on VGGFace dataset.

Methods	FID (\downarrow)	IS (\uparrow)
FIGR [4]	154.21	5.19
GMN [5]	201.12	6.38
DAGAN [7]	121.43	4.12
Ours	108.56	8.32

For our MatchingGAN, we train the model based on seen categories. Then, we randomly select $K = 3$ images from each unseen category, after which these conditional images and a random vector are fed into the trained model to generate a new image for this unseen category. We can repeat the above procedure to generate adequate images for each unseen category. Similarly, we train GMN [5] and FIGR [4] in 1-way 3-shot setting based on seen categories, and use trained model to generate images for unseen categories. Distinctive from the above methods, DAGAN [7] is conditioned on a single image, but we can still generate adequate images for unseen categories by using one conditional image each time.

We generate 128 images for each unseen category using each method, based on which FID and IS are calculated. The results of different methods are reported in Table 1, from which we observe that the images generated by our MatchingGAN achieve the highest IS and lowest FID, which demonstrates that our model could produce more diverse and realistic images compared with baseline methods.

For visualization comparison, we show some example images generated by our MatchingGAN on three datasets in Fig. 2. We also show the images generated by DAGAN for comparison, which is a competitive baseline as demonstrated in Table 1. It can be seen that our method can produce more diverse images than DAGAN, because our method excels in fusing the information of more than one conditional image. More visualization results can be found in Supplementary.

Table 2. Accuracy(%) of different methods on different datasets in low-data setting.

Method	Dataset	Accuracy		
		5	10	15
Standard	Omniglot	66.22	81.87	83.31
FIGR [4]	Omniglot	69.23	83.12	84.89
GMN [5]	Omniglot	67.74	84.19	85.12
DAGAN [7]	Omniglot	88.81	89.32	95.38
Ours	Omniglot	89.03	90.92	96.29
Standard	EMNIST	83.64	88.64	91.14
FIGR [4]	EMNIST	85.91	90.08	92.18
GMN [5]	EMNIST	84.56	91.21	92.09
DAGAN [7]	EMNIST	87.45	94.18	95.58
Ours	EMNIST	91.75	95.91	96.29
Standard	VGGFace	8.82	20.29	39.12
FIGR [4]	VGGFace	6.12	18.84	32.13
GMN [5]	VGGFace	5.23	15.61	35.48
DAGAN [7]	VGGFace	19.23	35.12	44.36
Ours	VGGFace	21.12	40.95	50.12

4.3. Low-data Classification

To further evaluate the quality of generated images, we use generated images to help downstream classification tasks in low-data setting in this section and few-shot setting in Section 4.4. For low-data classification on unseen categories, we randomly select a few (*e.g.*, 5, 10, 15) training images per unseen category while the remaining images in each unseen category are test images. Note that we have training and testing phases for classification, which are different from the training and testing phases of our MatchingGAN. We use ResNet18 [31] pretrained on seen categories as backbone network, train the classifier based on the training images of unseen categories, and finally predict the test images of unseen categories. This setting is referred to as “standard” in Table 2.

Then, we attempt to use generated images to augment the training set of unseen categories. For each few-shot generation method, we generate 512 images for each unseen category based on the training set of unseen categories. Then, the ResNet18 classifier is trained on the augmented training set (original training set and generated images) and applied to the test set of unseen categories. The results of different methods are listed in Table 2. On Omniglot and EMNIST datasets, all methods outperform “standard”, which demonstrates the benefit of augmented training set. On VGGFace dataset, our MatchingGAN and DAGAN [7] outperform “standard”, while GMN and FIGR underperform “standard”. One possible explanation is that the images generated by GMN and FIGR on the more challenging VGGFace dataset are of low quality and mislead the training of ResNet18. It can also be seen that our proposed MatchingGAN achieves significant improvement over baseline methods, which corroborates the

Table 3. Accuracy(%) of different methods on different datasets in few-shot setting.

Methods	Dataset	5-way 5-shot	10-way 5-shot
MatchingNets [9]	Omniglot	98.70	98.91
MAML [17]	Omniglot	99.90	99.13
RelationNets [14]	Omniglot	99.80	99.22
MTL [27]	Omniglot	99.85	99.35
DN4 [28]	Omniglot	99.83	99.29
Ours	Omniglot	99.93	99.42
MatchingNets [9]	VGGFace	60.01	48.67
MAML [17]	VGGFace	61.09	47.89
RelationNets [14]	VGGFace	60.93	49.12
MTL [27]	VGGFace	63.67	51.94
DN4 [28]	VGGFace	62.89	51.58
Ours	VGGFace	65.12	53.21

effectiveness of combining matching procedure with adversarial learning.

4.4. Few-shot Classification

Following the N -way C -shot setting in few-shot classification [9, 14], we create episodes and report the averaged accuracy over multiple episodes on each dataset. In each episode, we first randomly select N unseen categories, and then randomly select C images from each unseen category as training set and the remaining images are used as test set. Each episode is similar to the low-data setting in Section 4.3. Again, we use ResNet18 as the classifier and generate 512 images for each unseen category to augment the training set.

We compare our MatchingGAN with few-shot classification methods, including representative methods MatchingNets [9], RelationNets [14], MAML [17] as well as state-of-the-art methods MTL [27], DN4 [28]. Note that the above few-shot classification methods do not generate images to augment the training set of unseen categories. Instead, we strictly follow their original training procedure based on seen categories and fine-tuning procedure based on the training set of unseen categories if necessary.

We conduct experiments in 5-way/10-way 5-shot setting on Omniglot and VGGFace datasets, and report the averaged results over 10 episodes on each dataset. From Table 3, we can observe that our MatchingGAN achieves better results than few-shot classification methods, which shows the power of augmented images generated by our model.

4.5. Ablation Studies

We analyze the impact of hyper-parameters (*i.e.*, λ_r , λ_m , K) in our method and investigate different design choices like skip connection. We also remove the matching procedure and use random interpolation coefficients to prove the necessity

of matching procedure. Due to space limitation, we leave the detailed experimental results to Supplementary.

5. CONCLUSION

In this paper, we have proposed a novel few-shot generation method MatchingGAN by combining matching procedure with adversarial learning. Comprehensive generation and classification experiments on three datasets have demonstrated that our MatchingGAN can generate more diverse and realistic images than existing methods.

F. DETAILS OF NETWORK ARCHITECTURE

Matching Generator In our matching generator, E_ψ and G form an auto-encoder, which is a combination of UNet [22] and ResNet [31]. Specifically, the auto-encoder has 8 blocks (4 blocks for encoder and 4 blocks for decoder), in which each block contains 4 composite layers (leaky ReLU and batch normalization followed by one downscaling or upscaling layer). Downscaling layers (in blocks 1 – 4) are convolutional layers with stride 2 followed by leaky ReLU, batch normalization, and dropout. Upscaling layers (in blocks 5 – 8) are stride 1/2 replicators followed by a convolutional layer, leaky ReLU, batch normalization, and dropout. The first 2 blocks of encoder and the last 2 blocks of decoder have 64 convolutional filters, while the last 2 blocks of the encoder and the first 2 blocks of the decoder have 128 convolutional filters. Skip connections are added between the last two blocks of encoder and the first two blocks of decoder.

For matching procedure, E_ϕ has the same structure and shared model parameters with E_ψ . E_z only has a fully connected (fc) layer with d outputs, where d is the same dimension as the output of encoder E_ϕ .

Matching Discriminator The network structure of our discriminator is similar to that in [13]. The discriminator consists of one convolutional layer followed by five blocks with increasing numbers of channels. The structure of each block is as follows: ResBlk- $k \rightarrow$ ResBlk- $k \rightarrow$ AvePool2x2, where ResBlk- k is a ReLU first residual block [32] with the number of channels k set as 64, 128, 256, 512, 1024 in five blocks. We use one fully connected (fc) layer with 1 output following AvePool layer to obtain the discriminator score. The classifier shares the feature extractor with discriminator and only replaces the last fc layer with another fc layer with C^s outputs with C^s being the number of seen categories. To obtain the features for feature matching loss, we remove the last fc layer from discriminator to extract the discriminative features of conditional image \mathcal{X}_S and generated image $\tilde{\mathbf{x}}$.

G. DETAILS OF PERFORMANCE METRICS

Inception Scores We use the Inception Score (IS) [29], which is widely used for evaluating the quality of generated images.

Table 4. Accuracy(%) of low-data (10-sample) classification augmented by our MatchingGAN with different K_1 and K_2 on EMNIST dataset.

	$K_1 = 3$	$K_1 = 5$	$K_1 = 7$	$K_1 = 9$
$K_2 = 3$	95.91	95.72	94.89	93.96
$K_2 = 5$	94.01	96.11	95.79	95.08
$K_2 = 7$	92.89	93.89	96.92	96.16
$K_2 = 9$	88.42	90.12	92.21	97.11

Let $p(y|\tilde{\mathbf{x}})$ be the posterior distribution of generated image $\tilde{\mathbf{x}}$ over unseen categories. The inception score is given by:

$$IS = \exp \left(\mathbb{E}_{\tilde{\mathbf{x}} \sim \iota(\tilde{\mathbf{x}})} [\text{KL}(\iota(y|\tilde{\mathbf{x}}) | \iota(y))] \right) \quad (9)$$

where $\iota(y) = \int_{\tilde{\mathbf{x}}} \iota(y|\tilde{\mathbf{x}}) d\tilde{\mathbf{x}}$. It is argued that the inception score is positively correlated with visual quality of generated images.

Frchet Inception Distance The Frchet Inception Distance (FID) [30] is designed for measuring similarities between two sets of images. We remove the last average pooling layer of the ImageNet-pretrained Inception-V3 [33] model as the feature extractor. Then, we compute FID between the generated images and the real images from the unseen categories.

H. MORE EXAMPLE IMAGES GENERATED BY OUR MATCHINGGAN

We show more example images generated by our MatchingGAN ($K = 3$) on Omniglot, EMNIST, and VGGFace datasets in Fig. 3, Fig. 4, and Fig. 5 respectively. Besides, we additionally conduct experiments on FIGR [4] dataset, which is not used in our main paper. The generated images on FIGR dataset are shown in Fig. 6. On all four datasets, our MatchingGAN can generate diverse and plausible images based on a few conditional images.

I. ABLATION STUDIES

The number of conditional images To analyze the impact of the number of conditional images, we train MatchingGAN with K_1 conditional images based on seen categories, and generate new images for unseen categories with K_2 conditional images. By default, we set $K = K_1 = K_2 = 3$ in our experiments. We evaluate the quality of generated images using different K_1 and K_2 in low-data (*i.e.*, 10-sample) classification, which is the same as Section 4.3 in the main paper. By taking EMNIST dataset as an example, we report the results in Table 4 by varying K_1 and K_2 in the range of [3, 5, 7, 9]. From Table 4, we can observe that when $K_2 = K_1$, our MatchingGAN can generally achieve good performance and the performance increases as K increases. Besides, we observe that given a fixed K_2 , when $K_1 > K_2$, the performance is not degraded a lot compared with $K_2 = K_1$. However, given a fixed K_2 , when $K_1 < K_2$, the performance is

significantly degraded. We conjecture that if our MatchingGAN is trained with K_1 conditional images, it cannot generalize well to fuse the information of more conditional images ($K_2 > K_1$) in the testing phase.

Hyper-parameter analysis In our MatchingGAN, we add a hyper-parameter λ_r before the weighted reconstruction loss \mathcal{L}_1 and a hyper-parameter λ_m before the feature matching loss \mathcal{L}_m . We investigate the impact of hyper-parameters λ_r and λ_m on VGGFace dataset, by varying λ_r (*resp.*, λ_m) in the range of $[0.01, 0.1, 1, 10]$ (*resp.*, $[0, 1]$). We evaluate the quality of generated images from two perspectives. On one hand, we compute the Inception Score (IS) and Frchet Inception Distance (FID) of generated images as done in Section 4.2 in the main paper. On the other hand, we report the accuracy of low-data (10-sample) classification augmented with generated images as described in Section 4.3 in the main paper. The results are reported in Table 5, which shows that larger λ_r leads to lower FID and higher IS at the cost of classification performance. $\lambda_r = 0.1$ achieves a good trade-off, so we use $\lambda_r = 0.1$ as default value in our experiments. Another observation is that after removing the feature matching loss by setting $\lambda_m = 0$, IS, FID, and classification accuracy become significantly worse, which indicates the benefit of feature matching loss.

Random interpolation coefficient In our MatchingGAN, we employ matching procedure to learn reasonable interpolation coefficients, which are used to fuse the features of conditional images. A naive alternative to the matching procedure is to sample normalized random vector $[a_1, a_2, \dots, a_K]$ with $a_i \geq 0$ and $\sum_i a_i = 1$ from uniform distribution as interpolation coefficients. In this way, we can use random interpolation coefficients in both training and testing phase, so that the matching procedure could be discarded. In particular, E_z and E_ϕ will not participate in generator training. When using the trained generator to generate new images for unseen categories, we can also randomly sample interpolation coefficients without relying on E_z and E_ϕ .

By taking VGGFace dataset as an example, we compare MatchingGAN with matching procedure with the one without matching procedure (use random interpolation coefficients) based on three evaluation metrics (*i.e.*, FID, IS, and low-data classification accuracy). The results are listed in Table 5, which demonstrate that matching procedure is capable of learning more reasonable interpolation coefficients than random interpolation coefficients, leading to better generated images.

Network design choices To investigate surrogate choices of network design, we again take VGGFace dataset as an example and utilize three evaluation metrics (*i.e.*, FID, IS, and low-data classification accuracy) for comparison. In our MatchingGAN, the encoder E_ϕ has the same network structure as E_ψ with shared model parameters. Alternatively, we can learn different model parameters for E_ϕ and E_ψ separately. Table 5 records the results of MatchingGAN in these two

Table 5. Analyses of hyper-parameters and different network design choices on VGGFace dataset.

setting	accuracy	FID (\downarrow)	IS (\uparrow)
$\lambda_r = 0.01$	35.62	112.16	7.89
$\lambda_r = 0.1$	40.95	108.56	8.32
$\lambda_r = 1$	33.89	107.16	9.17
$\lambda_r = 10$	30.12	106.12	11.04
$\lambda_m = 1$	40.95	108.56	8.32
$\lambda_m = 0$	28.98	111.4	7.56
matching coefficient	40.95	108.56	8.32
random coefficient	38.12	110.98	7.92
shared encoder	40.95	108.56	8.32
different encoder	40.98	107.98	8.56
1 connection	38.67	113.21	7.09
2 connection	40.95	108.56	8.32
3 connection	34.12	106.12	9.14

different cases. We observe that although introducing more model parameters, learning two encoders separately does not notably improve the performance of MatchingGAN.

Besides, in our matching generator, the number of skip connections between encoder E_ψ and decoder G also affects the quality of fused features and generated images. We utilize two connection blocks by default in our experiments. Here, we further explore the effect of using different numbers of skip connections. We report the results using 1, 2, 3 skip connections in Table 5. For 1 skip connection, we only keep the skip connection between the last block in encoder E_ψ and the first block in decoder G . For 3 skip connections, we add another connection between the second block in encoder E_ψ and the third block in decoder G . According to Table 5, we can see that using more skip connections could improve the realism of generated images (lower FID and higher IS). Another observation is that 3 skip connections compromise the low-data classification performance, because the generated images become closer to conditional images and lacking of diversity based on our experimental observation. We conjecture that it would be better to fuse multi-level information with an appropriate number of skip connections, so we opt for two skip connections in our matching generator.

J. REFERENCES

- [1] Jianmin Bao, Dong Chen, Fang Wen, Houqiang Li, and Gang Hua, "CVAE-GAN: fine-grained image generation through asymmetric training," in *ICCV*, 2017.
- [2] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio, "Generative adversarial nets," in *NeurIPS*, 2014.
- [3] Martin Arjovsky, Soumith Chintala, and Léon Bottou, "Wasserstein generative adversarial networks," in *ICML*, 2017.



Fig. 3. Images generated from trained MatchingGAN with $K = 3$ on Omniglot Dataset. The real conditional images are the left three columns.

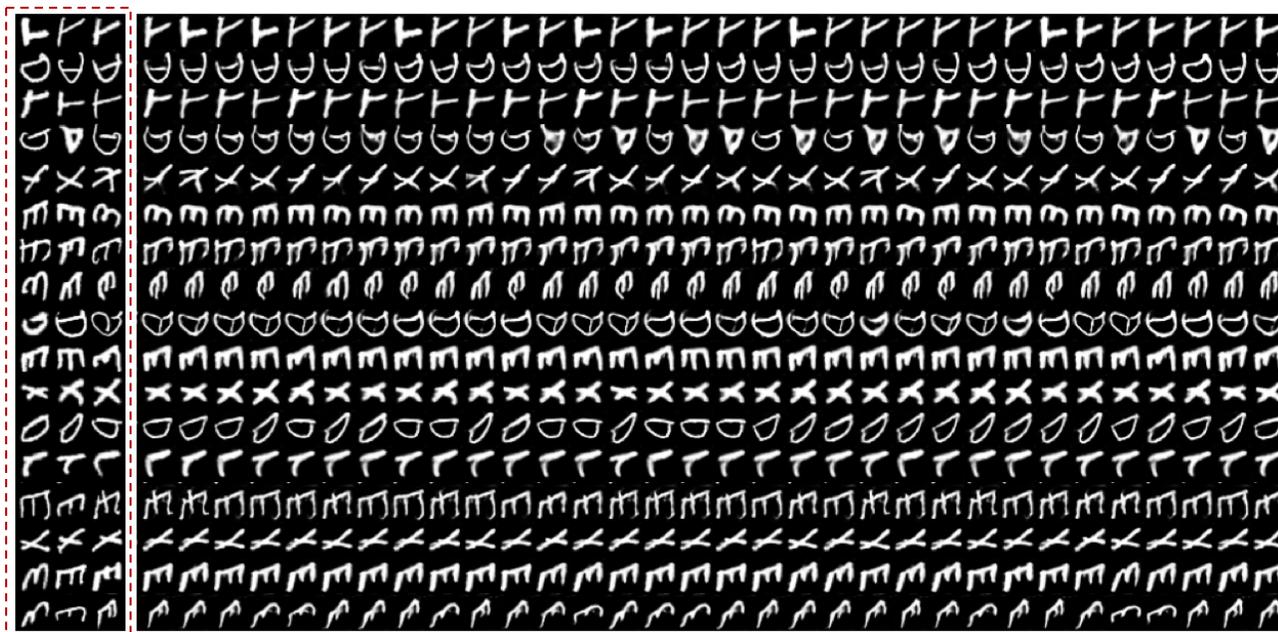


Fig. 4. Images generated from trained MatchingGAN with $K = 3$ on EMNIST Dataset. The real conditional images are the left three columns.

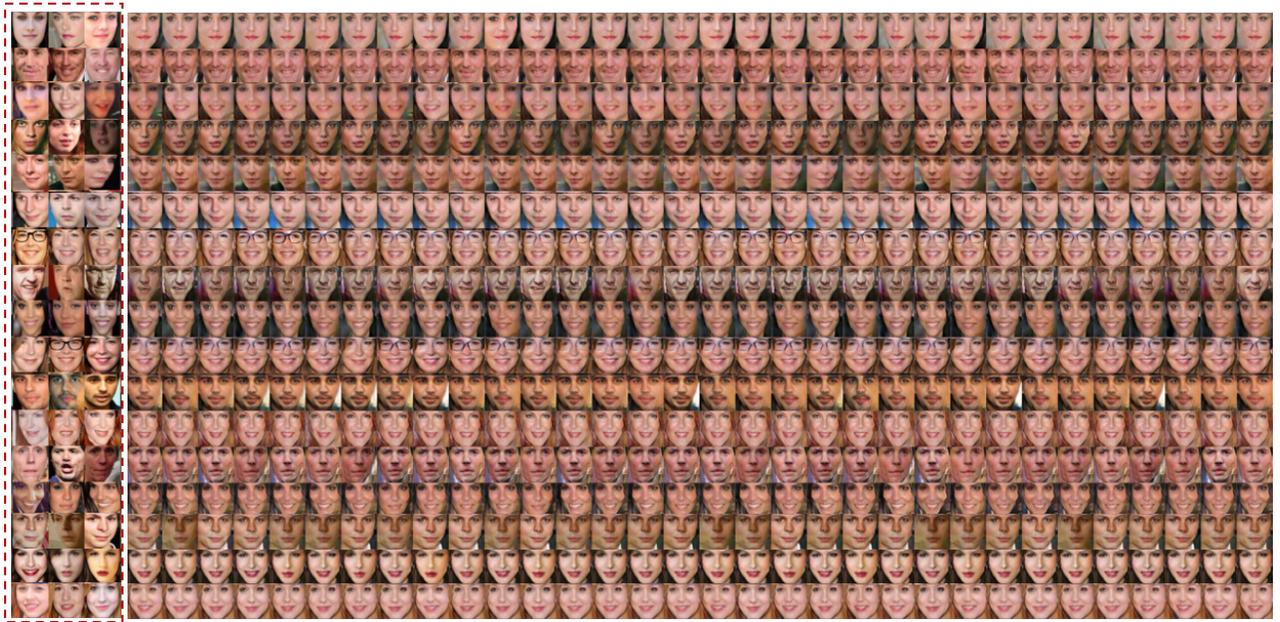


Fig. 5. Images generated from trained MatchingGAN with $K = 3$ on VGGFace Dataset. The real conditional images are the left three columns.

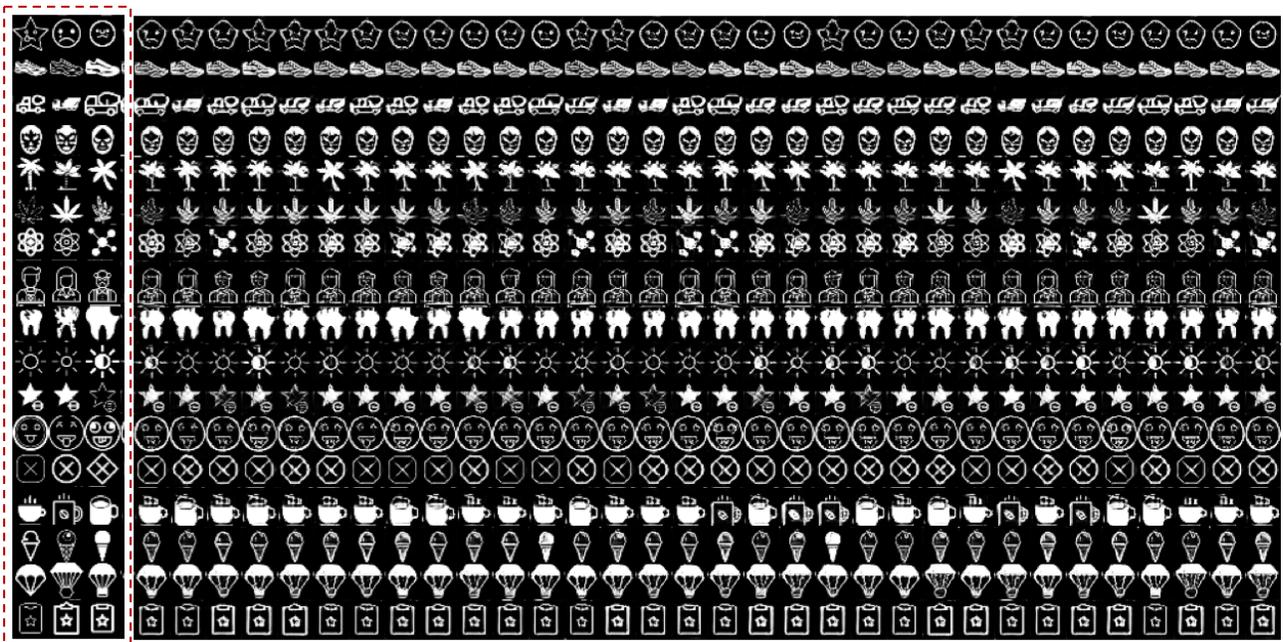


Fig. 6. Images generated from trained MatchingGAN with $K = 3$ on FIGR Dataset. The real conditional images are the left three columns.

- [4] Louis Clouâtre and Marc Demers, “Figr: Few-shot image generation with reptile,” *arXiv preprint arXiv:1901.02199*, 2019.
- [5] Sergey Bartunov and Dmitry Vetrov, “Few-shot generative modelling with generative matching networks,” in *ICAIIS*, 2018.
- [6] Eli Schwartz, Leonid Karlinsky, Joseph Shtok, Sivan Harary, Mattias Marder, Abhishek Kumar, Rogerio Feris, Raja Giryes, and Alex Bronstein, “Delta-encoder: an effective sample synthesis method for few-shot object recognition,” in *NeurIPS*, 2018.
- [7] Antreas Antoniou, Amos Storkey, and Harrison Edwards, “Data augmentation generative adversarial networks,” *arXiv preprint arXiv:1711.04340*, 2017.
- [8] Alex Nichol, Joshua Achiam, and John Schulman, “On first-order meta-learning algorithms,” *arXiv preprint arXiv:1803.02999*, 2018.
- [9] Oriol Vinyals, Charles Blundell, Timothy Lillicrap, Daan Wierstra, et al., “Matching networks for one shot learning,” in *NeurIPS*, 2016.
- [10] Takeru Miyato, Toshiki Kataoka, Masanori Koyama, and Yuichi Yoshida, “Spectral normalization for generative adversarial networks,” *arXiv preprint arXiv:1802.05957*, 2018.
- [11] Takeru Miyato and Masanori Koyama, “cGANs with projection discriminator,” *arXiv preprint arXiv:1802.05637*, 2018.
- [12] Yongyi Lu, Shangzhe Wu, Yu-Wing Tai, and Chi-Keung Tang, “Image generation from sketch constraint using contextual GAN,” in *ECCV*, 2018.
- [13] Ming-Yu Liu, Xun Huang, Arun Mallya, Tero Karras, Timo Aila, Jaakko Lehtinen, and Jan Kautz, “Few-shot unsupervised image-to-image translation,” *arXiv preprint arXiv:1905.01723*, 2019.
- [14] Flood Sung, Yongxin Yang, Li Zhang, Tao Xiang, Philip HS Torr, and Timothy M Hospedales, “Learning to compare: Relation network for few-shot learning,” in *CVPR*, 2018.
- [15] Wenbo Zheng, Lan Yan, Chao Gou, Wenwen Zhang, and Fei-Yue Wang, “A relation network embedded with prior features for few-shot caricature recognition,” in *ICME*, 2019.
- [16] Adam Santoro, Sergey Bartunov, Matthew Botvinick, Daan Wierstra, and Timothy Lillicrap, “One-shot learning with memory-augmented neural networks,” *arXiv preprint arXiv:1605.06065*, 2016.
- [17] Chelsea Finn, Pieter Abbeel, and Sergey Levine, “Model-agnostic meta-learning for fast adaptation of deep networks,” in *ICML*, 2017.
- [18] Xuefeng Du, Dexing Zhong, and Pengna Li, “Low-shot palm-print recognition based on meta-siamese network,” in *ICME*, 2019.
- [19] Alex Krizhevsky, I. Sutskever, and G. Hinton, “Imagenet classification with deep convolutional neural networks,” *NeurIPS*, 2012.
- [20] Bharath Hariharan and Ross Girshick, “Low-shot visual recognition by shrinking and hallucinating features,” in *ICCV*, 2017.
- [21] Yu-Xiong Wang, Ross Girshick, Martial Hebert, and Bharath Hariharan, “Low-shot learning from imaginary data,” in *CVPR*, 2018.
- [22] Olaf Ronneberger, Philipp Fischer, and Thomas Brox, “U-net: Convolutional networks for biomedical image segmentation,” in *MICCAI*. Springer, 2015.
- [23] Augustus Odena, Christopher Olah, and Jonathon Shlens, “Conditional image synthesis with auxiliary classifier GANs,” in *ICML*, 2017.
- [24] Brenden M. Lake, Ruslan Salakhutdinov, and Joshua B. Tenenbaum, “One-shot learning by inverting a compositional causal process,” *NeurIPS*, 2015.
- [25] Gregory Cohen, Saeed Afshar, Jonathan Tapson, and André van Schaik, “EMNIST: an extension of MNIST to handwritten letters,” *arXiv preprint arXiv:1702.05373*, 2017.
- [26] Qiong Cao, Li Shen, Weidi Xie, Omkar M Parkhi, and Andrew Zisserman, “Vggface2: A dataset for recognising faces across pose and age,” in *FG*, 2018.
- [27] Qianru Sun, Yaoyao Liu, Tat-Seng Chua, and Bernt Schiele, “Meta-transfer learning for few-shot learning,” in *CVPR*, 2019.
- [28] Wenbin Li, Lei Wang, Jinglin Xu, Jing Huo, Yang Gao, and Jiebo Luo, “Revisiting local descriptor based image-to-class measure for few-shot learning,” in *CVPR*, 2019.
- [29] Qiantong Xu, Gao Huang, Yang Yuan, Chuan Guo, Yu Sun, Felix Wu, and Kilian Weinberger, “An empirical study on evaluation metrics of generative adversarial networks,” *arXiv preprint arXiv:1806.07755*, 2018.
- [30] Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter, “GANs trained by a two time-scale update rule converge to a local nash equilibrium,” in *NeurIPS*, 2017.
- [31] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun, “Deep residual learning for image recognition,” in *CVPR*, 2016.
- [32] Lars Mescheder, Andreas Geiger, and Sebastian Nowozin, “Which training methods for GANs do actually converge?,” *arXiv preprint arXiv:1801.04406*, 2018.
- [33] Christian Szegedy, Vincent Vanhoucke, Sergey Ioffe, Jon Shlens, and Zbigniew Wojna, “Rethinking the inception architecture for computer vision,” in *CVPR*, 2016.