

Dense-View GEIs Set: View Space Covering for Gait Recognition based on Dense-View GAN

Rijun Liao¹, Weizhi An², Shiqi Yu³, Zhu Li¹, Yongzhen Huang^{4,5}

¹Department of Computer Science and Electrical Engineering, University of Missouri-Kansas City, USA.

²Department of Computer Science and Engineering, University of Texas at Arlington, USA

³Department of Computer Science and Engineering, Southern University of Science and Technology, China

⁴National Laboratory of Pattern Recognition, Institute of Automation, Chinese Academy of Sciences, China

⁵Watrix technology limited co. ltd, China

rijun.liao@mail.umkc.edu, weizhi.an@mavs.uta.edu, yusq@sustech.edu.cn,
lizhu@umkc.edu, yongzhen.huang@nlpr.ia.ac.cn

Abstract

Gait recognition has proven to be effective for long-distance human recognition. But view variance of gait features would change human appearance greatly and reduce its performance. Most existing gait datasets usually collect data with a dozen different angles, or even more few. Limited view angles would prevent learning better view invariant feature. It can further improve robustness of gait recognition if we collect data with various angles at 1° interval. But it is time consuming and labor consuming to collect this kind of dataset. In this paper, we, therefore, introduce a Dense-View GEIs Set (DV-GEIs) to deal with the challenge of limited view angles. This set can cover the whole view space, view angle from 0° to 180° with 1° interval. In addition, Dense-View GAN (DV-GAN) is proposed to synthesize this dense view set. DV-GAN consists of Generator, Discriminator and Monitor, where Monitor is designed to preserve human identification and view information. The proposed method is evaluated on the CASIA-B and OU-ISIR dataset. The experimental results show that DV-GEIs synthesized by DV-GAN is an effective way to learn better view invariant feature. We believe the idea of dense view generated samples will further improve the development of gait recognition.

1. Introduction

Gait is one kind of popular biometric features for human identification. Compared with other features like face, iris, palmprint and fingerprint, gait provides a unique possibil-

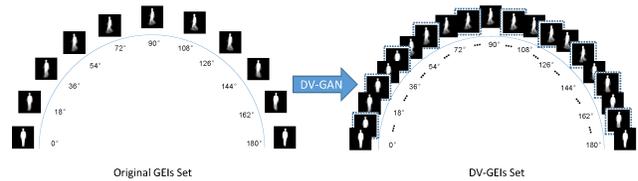


Figure 1. Dense-View GEIs Set (DV-GEIs): view space covering for learning better view invariant feature, view angle range from 0° to 180° with 1° interval. DV-GAN is proposed to synthesize realistic samples with various view angle conditions. (Sample images from CASIA-B dataset [24])

ity to identify a subject at a long distance without people's cooperation. Therefore, it has great potential application in catching criminals, video surveillance and social security.

However, gait recognition is still challenging in the real application. This is because there are many variations would reduce its performance. View is one of gait challenges because we can not control the walking direction of people and view changing that lead the human body shape to change greatly. In this paper, we propose DV-GEIs to further reduce the influence of view variation.

In order to solve the cross-view challenge, some researchers focus on view transformation model (VTM) which can transform gait features from one view to another view. Examples are FD-VTM [15], RSVD-VTM [11] and RPCA-VTM [26]. But most VTM methods need to know angles of probe and gallery before extracting gait features. This means that each view needs one model which has some challenges in the real application. In the next few years, SPAE [25], GaitGAN [22] and GaitGANv2 [23] are proposed that can transform any view gait into the side view gait by using only an uniform model. However, the side

view transform strategy will collapse when the view variance is large.

With the development of deep learning, some works [13, 1, 14, 2] only use several human pose coordinates as input gait features which is robust to human appearance. However, its performance still need to be improved because human pose coordinates have not enough information compared with human appearance. In addition, Wu *et al.* [21] and Chao *et al.* [3] directly take a sequence of human silhouettes as input data rather than using the hard-crafted gait features and achieve high performance. The price of these methods is high computational cost.

The above methods have made a great contribution to the development of gait recognition. We think that it can further improve its robustness to view variation if we can collect data with more view angle. Because most existing datasets not cover all kinds of view condition, and deep learning depends heavily on big data. CASIA-B dataset [24] has 11 views with 18° view angle interval, OU-MVPL [19] has 14 views and OU-ISIR [9] only has four views. It has negative influence to learn better invariant feature if the number of view angle in the dataset is limited. We can collect data with more view angles manually. But it is time consuming and labor consuming to collect dense view angle data.

Recently, some researchers use synthesized samples based on GAN to improve original performance. Chen *et al.* [4] used GAN to generate noise samples and use them in image blind denoising. Qian *et al.* [17] combined human pose and GAN to synthesize human image in specific pose for person re-identification. Those methods can greatly improve its original performance. But this idea has not been achieved in gait recognition task, because it needs to find a suitable solution to synthesize samples with different conditions. One popular based on GAN work, GaitGAN [22], has used GAN to transform any view GEI into the side view GEI, which effectively improve gait robustness. In our work, we extend the development of GAN in gait recognition by generating samples with more views, rather than view transformation.

Unlike above view transformation methods, we generate gait features with dense views to improve recognition rate on the cross-view condition. We are inspired by the idea of synthesized samples and GaitGAN [22], and proposed a novel generation model DV-GAN to synthesize DV-GEIs set. GEI [6] is employed to be the gait features in the proposed method same as GaitGAN [22] and SPAE [25], because its robustness to noise and its efficiency in computation. Our method in this paper has the following contributions:

- We introduce a novel Dense-View GEIs Set (DV-GEIs) to solve the challenge of limited view angles on the existing gait dataset. Most existing datasets usually capture samples with several or dozens view at large in-

terval, while the view angle of DV-GEIs could cover whole view space to make up for small number of view angles, range from 0° to 180° with 1° interval, as shown in Figure 1.

- A novel GEI generation model Dense-View GAN (DV-GAN) is proposed to generate realistic GEIs with various view angle conditions. Compared with traditional GAN [5] which has generator and discriminator, DV-GAN includes additional monitor which can maintain human identification and view information very well.
- We have performed several experiments on CASIA-B [24] and OU-ISIR [9] dataset. Experimental results shows that dense view samples synthesized by DV-GAN can further improve robustness to view variation compared with original dataset.

2. Method

2.1. Dense-View GEIs Set (DV-GEIs)

We denote the views of GEI as p, q which are correspondent with the input GEI x_p, x_q . Our aim is to synthesize various views GEIs, from p angle to q angle, as shown in Figure 2. The synthesized GEIs are created according to the following equation:

$$x' = \{G_D(z) | z = \alpha z_p + (1 - \alpha)z_q\} \quad (1)$$

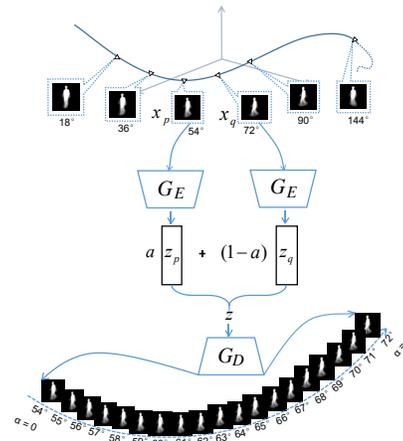


Figure 2. The latent space z_p, z_q are encoded by encoder G_E with two input GEIs x_p, x_q . We synthesize various views GEIs from p angle to q angle by decoding the latent space z , where $z = \alpha z_p + (1 - \alpha)z_q, \alpha \in [0, 1]$.

where $z_p = G_E(x_p), z_q = G_E(x_q)$, latent space z_p, z_q are encoded by encoder G_E which keep the characteristic of gait attribute. The interpolation is defined by linear transformation $z = \alpha z_p + (1 - \alpha)z_q$, where $\alpha \in [0, 1]$, and then z is fed to decoder G_D to generate new GEIs.

The idea of DV-GEIs is inspired by the idea of Hou *et al.* [7] which can generate a series of different view angle human faces from *left face* to *right face* by linear transformation $z = \alpha z_p + (1 - \alpha)z_q$ in latent space. In [7], authors model the face attribute distribution by training lots of faces with autoencoder, and produce latent vectors that can capture the semantic information of face expressions. In addition, Hou *et al.* [7] investigate the latent space and show that semantic relationship between different latent representations can be used in facial attribute prediction. In our work, we take advantage of latent space to deal with the challenge of samples with limited view angle in gait recognition task.

2.2. Dense-View GAN (DV-GAN)

We propose DV-GAN to generate realistic GEIs at various view angle conditions. With the rapid development of adversarial generative network (GAN) [5, 17], GAN can generate images with realistic details. Our DV-GAN model consists of three neural networks: generator G , monitor M and discriminator D as shown in Figure 3.

- **Generator:** Given an input GEI x , and a target GEI \hat{x} , where $x = \hat{x}$. The purpose of our generator is to reconstruct GEI and model gait attribute distribution in latent space. We borrow the pixels to pixels level idea [8] to reconstruct the GEI image, that is adding L_1 norm loss to ensure the output GEI $\hat{x} = G(z, x)$ is the same as input x , our generator loss is defined as:

$$\min_{E,G} L_{L1}(G(z), x) \quad (2)$$

U-Net structure (add skip connections between two layers) is employed in our generator, as U-Net [18] architecture allows low-level information to shortcut across the network and effectively improve the quality of the generated images. We divide U-Net into two parts, encoder $G_E\{e_1\}$ and decoder $G_D\{e_2, e_3, e_4, e_5, e_6, d_1, d_2, d_3, d_4, d_5\}$. The feature map of e_1 layer is defined as latent space z , because we use the U-Net structure network and it can not decode latent space if other layers' feature map as latent space z . For example, if we define e_2 output feature as latent space z , then $G_E\{e_1, e_2\}$ and $G_D\{e_3, e_4, e_5, e_6, d_1, d_2, d_3, d_4, d_5\}$ will as encoder and decoder respectively. We can do linear interpolation $z = \alpha G_E(x_p) + (1 - \alpha)G_E(x_q)$, but we can not decode latent space z , because the calculation of feature map of d_5 requires the feature map of e_1 in U-Net structure, while decoder G_D does not include e_1 layer.

- **Discriminator:** Our discriminator network D ensures that the generated GEIs are realistic. It takes a pair of a real image x and a synthesized image \hat{x} as input, and is trained to identify whether an image is real or not.

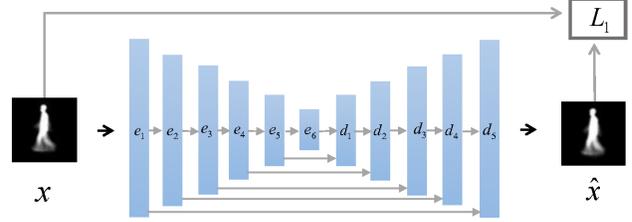


Figure 3. The U-net architecture generator of DV-GAN for modelling gait view space distribution.

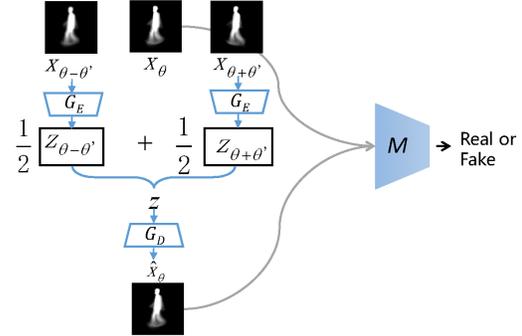


Figure 4. The structure of the real/fake monitor. Monitor will identify if the generate image \hat{x}_θ is same as the original image x_θ or not, to preserve human identification and view information.

If the input image is from a real image, the discriminator output value is 1, otherwise 0. The discriminator would ensure the generated GEI is more and more similar to the original image. The objective of discriminator is:

$$\min_G \max_D \mathbb{E}_{x,z} p_{data(x,z)} [\log D(x, x)] + \mathbb{E}_{x,z} p_{data(x,z)} [1 - \log D(x, G(z))] \quad (3)$$

- **Monitor:** Monitor is designed to preserve human identification and view information, as shown in Figure 4. This idea is inspired by identification-discriminator of GaitGAN [22], which takes a source image and a target image as input and is trained to identify whether the input pair is the same person. Our monitor not only can maintain the identification information, but also can preserve the view information. The monitor has three input images $x_{\theta-\theta'}$, x_θ and $x_{\theta+\theta'}$. In the training process, first, monitor would generate an image \hat{x}_θ by decode the latent feature z , where latent feature z is the mean value of encoded features ($z_{\theta-\theta'}$ and $z_{\theta+\theta'}$) of two input images ($x_{\theta-\theta'}$ and $x_{\theta+\theta'}$). And then the monitor would produce a scalar probability to indicate if the generate image \hat{x}_θ is same as the original image x_θ or not, so the identification and view information would be maintained in the training phase, the formula as follows:

$$\min_G \max_D \mathbb{E}_{x_\theta, z \sim p_{data}(x_\theta, z)} [\log D(x_\theta, x_\theta)] + \mathbb{E}_{x_\theta, z \sim p_{data}(x_\theta, z)} [1 - \log D(x_\theta, \hat{x}_\theta)]$$

$$\text{where } \hat{x}_\theta = G\left(\frac{1}{2}E(x_{\theta-\theta'}) + \frac{1}{2}E(x_{\theta+\theta'})\right) \quad (4)$$

3. Experiments and Analysis

3.1. Datasets

The proposed method is evaluated on CASIA-B dataset [24] with 11 views and OU-ISIR Large Population Dataset [9] with 4 views, respectively.

CASIA-B gait dataset [24] is one of the popular public gait databases and it was created by the Institute of Automation, Chinese Academy of Sciences in January 2005. It contains 124 subjects, and each subject has six sequences. This set is captured at 11 views, from 0° to 180° with 18° interval between two nearest views. The view angles are $\{0^\circ, 18^\circ, \dots, 180^\circ\}$. The left image in Figure 1 illustrates the samples involving 11 views from a normal walking subject. CASIA-B has three different conditions, normal walking (NM), walking with bag (BG) and walking with a coat (CL). Our proposed method focus on normal walking (NM) condition to solve the view variation.

In order to better evaluate our method, we perform another experiment on OU-ISIR Large Population Dataset [9] with only 4 views ($55^\circ, 65^\circ, 75^\circ$ and 85°). OU-ISIR is a very large dataset which contains 4007 subjects ranging from 1 to 94 years old. It includes two sequences under the normal walking conditions. It enables us to study the upper bound of gait recognition performance in a more statistically reliable manner.

3.2. Experimental Setting

For fair comparison with SPAE [25] and GaitGAN [22] later, our experimental setting is same as that of them. The training set contains the first 62 subjects under 6 normal sequences, and the test set contains the rest of subjects. In the test set, the gallery set contains the first 4 normal sequences and the probe set is consists of the rest 2 normal sequences, as shown in Table 1. We only synthesize GEIs in training set, not in the test set in our proposed method. Because we want to show that the model trained by dense view synthesized samples can improve performance in the test set compared with original dataset.

The setting of OU-ISIR [9] is similar to that of CASIA-B. In the experiment, we divide all the subjects into five sets randomly, and keep one set for testing and four sets for training to synthesize samples. In each test set, the first sequence is put into gallery set and the rest sequence is put into probe set.

Table 1. Experimental setting on CASIA-B dataset (NM: normal walking).

Training	Test	
	Gallery Set	Probe Set
ID: 001-062	ID: 063-124	ID: 063-124
NM01-NM06	NM01-NM04	NM05-NM06

Table 2. Implementation details of the Generator network

Layers	Number of filters	Filter size	Stride	Batch norm	Activation function
Conv.1	64	5×5	2	N	L-ReLU
Conv.2	128	5×5	2	Y	L-ReLU
Conv.3	256	5×5	2	Y	L-ReLU
Conv.4	512	5×5	2	Y	L-ReLU
Conv.5	512	5×5	2	Y	L-ReLU
Conv.6	512	5×5	2	Y	L-ReLU
Deconv.1	512	5×5	2	Y	ReLU
Deconv.2	512	5×5	2	Y	ReLU
Deconv.3	256	5×5	2	Y	ReLU
Deconv.4	128	5×5	2	Y	ReLU
Deconv.5	64	5×5	2	Y	ReLU
Deconv.6	64	5×5	2	N	Tanh

Table 3. Implementation details of the Discriminator network.

Layers	Number of filters	Filter size	Stride	Batch norm	Activation function
Conv.1	64	5×5	2	N	L-ReLU
Conv.2	128	5×5	1	Y	L-ReLU

Table 4. Implementation details of the Monitor network.

Layers	Number of filters	Filter size	Stride	Batch norm	Activation function
Conv.1	64	5×5	2	N	L-ReLU
Conv.2	128	5×5	1	Y	L-ReLU

3.3. Implementation Details of DV-GAN

Our structure of DV-GAN is inspired by the idea of Isola *et al.* [8]. Authors provide an open code, namely *pix2pix*, to solve the problem of image-to-image translation. This networks can learn the mapping from input image to output image, which is effective at synthesizing photos. The number of layers of our generator and discriminator is less than that of *pix2pix*, because our input size of image is 64×64 , while *pix2pix* is 256×256 . The implementation detail of generator and discriminator can be seen in Table 2 and Table 3. The output of discriminator is one dimensional, the convolution layer 2 is applied to map to a one dimensional output, followed by a sigmoid function.

In addition, we add additional monitor to preserve hu-

Table 5. Details implementation of the CNN.

Layers	Number of filters	Filter size	Stride	Activation function
Conv.1	32	3 × 3	1	P-ReLU
Conv.2	64	3 × 3	1	P-ReLU
Pooling.1	N	2 × 2	2	N
Conv.3	64	3 × 3	1	P-ReLU
Conv.4	64	3 × 3	1	P-ReLU
Eltwise.1	Sum operation between Pooling.1 and Conv.4			
Conv.5	128	3 × 3	1	P-ReLU
Pooling.2	N	2 × 2	2	N
Conv.6	128	3 × 3	1	P-ReLU
Conv.7	128	3 × 3	1	P-ReLU
Eltwise.2	Sum operation between Pooling.2 and Conv.7			
Conv.8	128	3 × 3	1	P-ReLU
Conv.9	128	3 × 3	1	P-ReLU
Eltwise.3	Sum operation between Eltwise.2 and Conv.9			
Conv.10	128	3 × 3	1	P-ReLU
FC.1	512	N	N	N

man identification and view information. The implementation detail of monitor (Table 4) is the same as that of discriminator, but their input data setting is different. The number of input image in discriminator is two, while monitor is three. In the experiment on CASIA-B, we set the $\theta' = 18^\circ$ in the Equation 4, where $\theta \in \{18^\circ, 36^\circ, 54^\circ, 72^\circ, 90^\circ, 108^\circ, 126^\circ, 144^\circ, 172^\circ\}$.

After we train the DV-GAN model, DV-GEIs set will be generated. We synthesize GEI from 0° to 180° with 1° interval by linear transformation $z = \alpha z_p + (1 - \alpha)z_q$ and decoder latent space $G_D(z)$. Follow the Equation 1, we set $\alpha \in \{\frac{1}{18}, \frac{2}{18}, \dots, \frac{17}{18}\}$, where the angle set of z_p and z_q is $\{(0^\circ, 18^\circ), (18^\circ, 36^\circ), \dots, (162^\circ, 180^\circ)\}$. So we can get the various view angle GEIs that do not exist in the original dataset. Finally, we combine the synthesized GEIs and original GEIs to form the DV-GEIs set, and fed into CNN to extract invariant feature.

3.4. Implementation Details of Feature Extraction

We use a simple CNN to extract view invariant feature from DV-GEIs set. We borrow the idea of CNN structure and multi-loss function of PoseGait [14] which can effectively extract gait dynamic and static information from human pose sequence. The network details can be seen in Table 5.

Our multi-loss function consists of center loss and softmax loss. The center loss [20] with the softmax loss jointly supervise the learning of our CNN. The multi-loss function is defined as Equation 5. The softmax loss is useful to pull apart different GEIs and it can enlarge the inter-class dispersion. The center loss would minimizing the intra-class variation and keep the features of different classes separa-

ble. In the CNN training process, each batch should calculate several centers by averaging the GEIs features of the corresponding label.

$$L = L_S + \gamma L_c$$

$$= - \sum_{i=1}^m \log \frac{e^{W_i^T \hat{x}_i + b_i}}{\sum_{j=1}^n e^{W_j^T \hat{x}_i + b_j}} + \frac{\gamma}{2} \sum_{i=1}^m \|\hat{x}_i - c_{li}\|_2^2 \quad (5)$$

where $\hat{x}_i \in \mathbb{R}^d$ is the i th GEI feature that belongs to the l_i th class. d , $W \in \mathbb{R}^{d \times n}$ and $b \in \mathbb{R}^d$ denote the feature dimension, last connected layer and bias term, respectively. $c_{li} \in \mathbb{R}^d$ is the l_i th class center of gait features. We set $\gamma = 0.008$ in the experiment.

3.5. Experimental Results on CASIA-B dataset

The experimental results on CASIA-B dataset are shown in Table 6. In this table, each row is correspondent to a GEI angle of the gallery, and each column is correspondent to the angle of the probe set. The recognition rate of the cross-view condition has 121 combinations.

In order to illustrate synthesized samples by our DV-GAN can make contributions to the improvement of gait recognition, we compare with another experiment OG-GEIs result. OG-GEIs model is trained by using original GEIs set. From Figure 5, we can see performance of DV-GEIs is better than OG-GEIs at many points. This shows that our dense view samples synthesized by DV-GAN is an effective way to learn better view invariant feature and enhance the robustness for gait recognition.

3.6. Visualization of Synthesized GEIs

In order to see the quality of generated image, we synthesize some view samples, as shown in Figure 6. In fact, in our above experiment, we generate the GEIs by using two adjacent GEIs with 18° interval in CASIA B dataset. However, the difference between the adjacent angle (18°) is hard to be distinguished in vision, and without ground truth GEIs in the original dataset. In order to visualize the transformation obviously and have ground truth for comparison, we use two GEIs (0° and 90°) with large 90° interval to generate some sample GEIs. That is, in the linear transformation $z = \alpha z_p + (1 - \alpha)z_q$, we set linear ratio $\alpha \in \{\frac{18}{90}, \frac{36}{90}, \frac{54}{90}, \frac{72}{90}\}$, and the angle set of z_p and z_q is $\{(0^\circ, 90^\circ)\}$.

To better show the contributions of DV-GAN, we compare our DV-GEIs with another type synthesized GEIs which direct view morphing by linear interpolation of two GEI images. That is, the second row GEIs are generated by equation $\hat{x} = \alpha x_p + (1 - \alpha)x_q$, where $\alpha \in \{\frac{18}{90}, \frac{36}{90}, \frac{54}{90}, \frac{72}{90}\}$. From Figure 6, we can see that synthesized GEIs by direct view morphing have obvious ghost, while synthesized GEIs

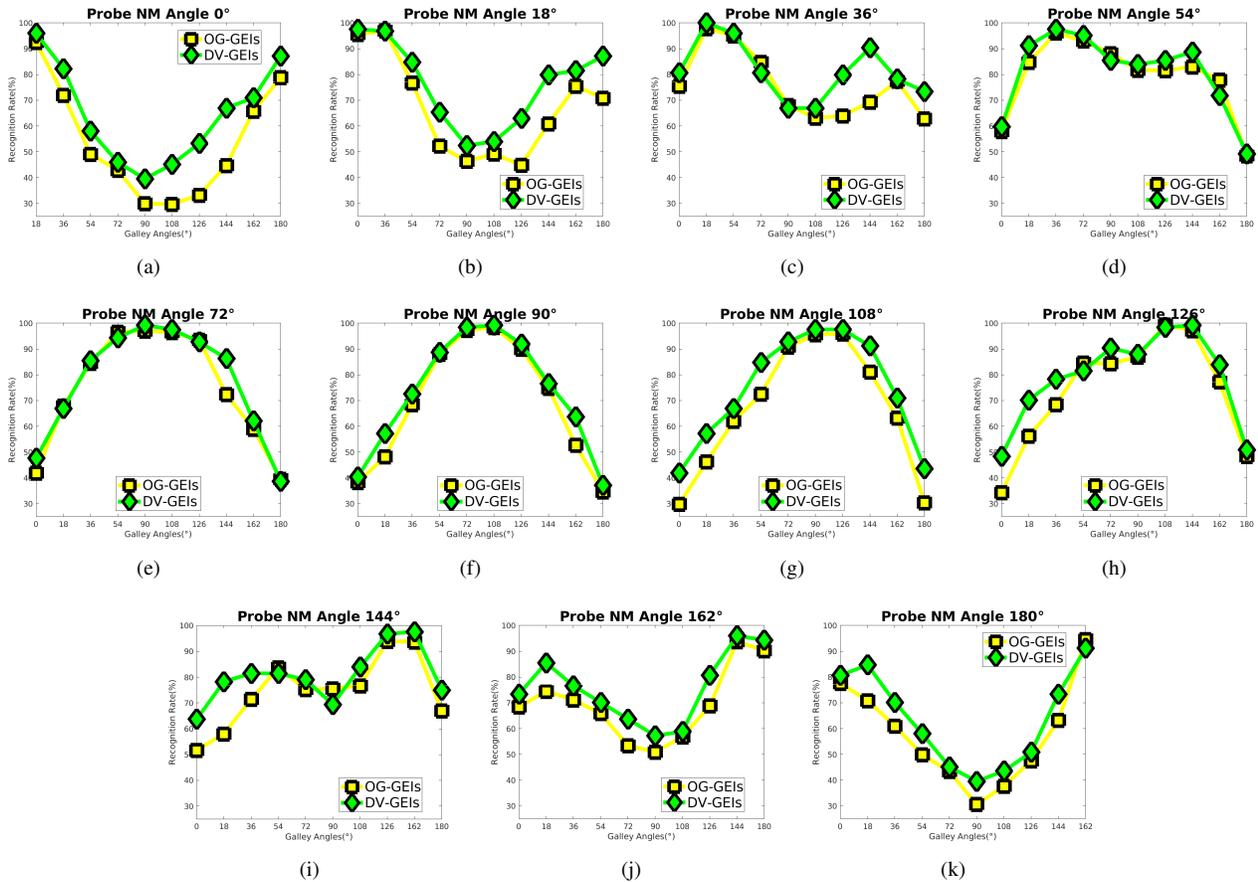


Figure 5. Comparison with OG-GEIs model that is trained by using original GEIs set, while DV-GEIs model is trained by using proposed dense view set. Each row represents a probe angle and each column represents different probe sequences in the test set. The comparison shows that samples with dense view synthesized by DV-GAN can further learn better view invariant feature.

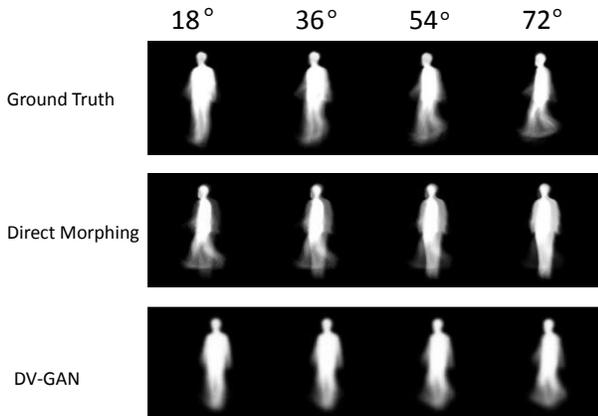


Figure 6. Visualization of synthesized GEIs. Second row GEIs are synthesized by linear interpolation of two GEI images. Third row GEIs are synthesized by proposed DV-GAN. Two types of synthesized GEIs are generated by 0° and 90° two original GEIs (see the visual demo in supplementary material).

by DV-GANs is very similar to ground truth although they are synthesized by two images with 90° interval. It will be more similar to the original image if using images with a smaller angle interval to synthesize. The comparison shows that our DV-GAN can synthesize realistic image with any angle condition, and our generated image can preserve human view information very well.

3.7. Comparison with VTM methods

In order to show the advantages of dense view samples, we compare with view transformation model (VTM) methods, including FD-VTM [16], RSVD-VTM [10], RPCA-VTM [27], R-VTM [12], SPAE [25], GaitGAN [22] and GaitGANv2 [23], as shown in Figure 7. Those methods are all trying to transform gait features from one view to another view, while our method is synthesizing more samples to cover the whole view space.

The probe angles selected are 54°, 90° and 126° in experiments of those methods. From Figure 7, we can see that the performance of proposed DV-GEIs method outper-

Table 6. Recognition rates of proposed method.(NM: normal walking).

		Probe set view (NM05,NM06)										
		0°	18°	36°	54°	72°	90°	108°	126°	144°	162°	180°
Gallery set view (NM01-NM04)	0°	100.0	97.58	80.65	59.68	47.58	40.32	41.94	48.39	63.71	73.39	80.65
	18°	95.97	100.0	100.0	91.13	66.94	57.26	57.26	70.16	78.23	85.48	84.68
	36°	82.26	96.77	98.39	97.58	85.48	72.58	66.94	78.23	81.45	76.61	70.16
	54°	58.06	84.68	95.97	97.58	94.35	88.71	84.68	81.45	81.45	70.16	58.06
	72°	45.97	65.32	80.65	95.16	99.19	98.39	92.74	90.32	79.03	63.71	45.16
	90°	39.52	52.42	66.94	85.48	99.19	99.19	97.58	87.90	69.35	57.26	39.52
	108°	45.16	54.03	66.94	83.87	97.58	99.19	99.19	98.39	83.87	58.87	43.55
	126°	53.23	62.90	79.84	85.48	92.74	91.94	97.58	97.58	96.77	80.65	50.81
	144°	66.94	79.84	90.32	88.71	86.29	76.61	91.13	99.19	99.19	95.97	73.39
	162°	70.97	81.45	78.23	71.77	62.10	63.71	70.97	83.87	97.58	98.39	91.13
180°	87.10	87.10	73.39	49.19	38.71	37.10	43.55	50.81	75.00	94.35	97.58	

Table 7. Comparison with based on GEI template methods on CASIA-B dataset at average accuracy(%). Excluding identical-view cases.

Training Subjects	Methods	Probe angle											
		0°	18°	36°	54°	72°	90°	108°	126°	144°	162°	180°	Mean
62	SPAE [25]	50.0	58.1	61.0	63.3	64.0	62.1	62.3	66.3	64.4	54.5	46.7	59.3
	GaitGAN [22]	41.9	53.5	63.0	64.5	63.1	58.1	61.7	65.7	62.7	54.1	40.6	57.2
	GaitGANv2 [23]	48.1	61.9	68.7	71.7	66.7	64.8	66.0	70.2	71.6	58.9	46.1	63.1
	DV-GEIs (Ours)	64.5	76.2	81.3	80.8	77.1	72.6	74.4	78.9	80.6	75.6	63.7	75.1
74	GaitSet-GEI [3]	-	-	-	-	-	-	-	-	-	-	-	80.4
	DV-GEIs (Ours)	71.0	86.4	91.4	89.6	80.4	80.1	82.5	90.1	90.4	85.3	70.5	83.4

forms that of others, especially when the angle difference between the gallery and the probe is large. This shows that dense view samples can deal with large viewpoint variation well. In addition, the proposed method can also improve the recognition rate obviously when the viewpoint variation is not large enough.

3.8. Comparison with GEI template methods

To further evaluate the proposed method, we compare with based on recent GEI template methods because our input data is also based on GEI, including SPAE [25], GaitGAN [22], GaitGANv2 [23], and GaitSet-GEI [3]. The comparison as shown in Table 7. Compared with those based on GEI template methods, our proposed method mean accuracy (75.1%) is much better than that of SPAE [25] (59.3%), GaitGAN [22] (57.2%) and GaitGANv2 [23] (63.1%). Those methods transform any GEI into the side GEI, while ours synthesize dense view samples to cover the whole view space. This show that gait view space covering can better handle with cross-view problem.

The method of GaitSet [3] can achieve a very high performance when it uses the human silhouette sequence as input feature. But its performance would be decreased dramatically when use GEI as input feature (drop from 95.0% to 80.4%). One reason for this is that human silhouette sequence has much rich information than GEI. Here, we not compare with those based on human silhouette sequence

method. Because our method is based on the GEI template, so it maybe unfair to compare with those based on human silhouette sequence method.

3.9. Experimental results on OU-ISIR dataset

OU-ISIR dataset [9] is also employed to evaluate the proposed method. In the training phase, the view angle of DV-GEIs on OU-ISIR is from 55° to 85° with 1° interval. This is because start angle if from 55°, and end with 85° on OU-ISIR dataset. In the test phase, probe angle and gallery angle are 55°, 65°, 75° and 85°, respectively. The result of experiments on OU-ISIR dataset is shown in Table 8. In this table, each row is correspondent to a angle of gallery set, and each column is correspondent to a angle of probe set. The recognition rate of the cross-view condition has 16 combinations.

We compare our results with OG-GEIs (trained by original GEIs set), DeepCNN [21] and GaitGANv2 [23], as shown in Table 9. In this table, each column is correspondent to the angle of the probe set. The recognition rate is by averaging different gallery angle, excluding identical view cases. From that table, we can see the performance of DV-GEIs is better than the baseline reported by the dataset authors [21, 23] when probe angle is 55° and 65°. In addition, the accuracy DV-GEIs outperforms that of OG-GEIs, which shows again that dense view samples synthesized by DV-GEIs can further improve the robustness to view variation.

Table 8. Experimental results on OU-ISIR dataset. Model is trained by using DV-GEIs dataset.

Probe angle	Gallery angle			
	55°	65°	75°	85°
55°	96.6	95.2	94.5	88.1
65°	97.3	96.8	94.8	93.2
75°	90.1	96.3	97.0	96.8
85°	91.1	96.3	96.4	96.5

Table 9. Comparison with other methods on OU-ISIR with average accuracy(%). Excluding identical view cases. OG-GEIs: trained by using original GEIs set. DV-GEIs: trained by using proposed dense view set.

Methods	Probe angle			
	55°	65°	75°	85°
DeepCNN [21]	91.6	92.3	92.4	94.8
GaitGANv2 [23]	91.9	95.0	94.4	94.6
OG-GEIs(Ours)	92.0	93.8	94.0	93.8
DV-GEIs(Ours)	92.6	95.1	94.4	94.6

4. Conclusions

In this paper, we introduce a novel Dense-View GEIs Set (DV-GEIs) to handle with the challenge of samples with limited view angles on existing gait datasets. View angle of DV-GEIs set can cover the whole view space, range from 0° to 180° with 1° interval. This set is synthesized by proposed DV-GAN, which consists of generator, discriminator and monitor. Monitor can preserve human identification and view information very well. The experimental result shows that samples with dense view can learn better view invariant feature compare with original dataset.

With the development of synthesized sample technology, we believe the idea of dense view samples synthesized by DV-GAN not only can enhance robustness to view variation, but also deal with other variations, like carrying and clothing condition. Eventually, it would further improve the development of gait recognition.

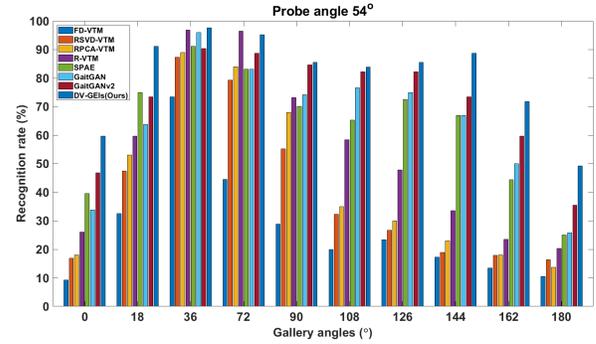
Acknowledgment

This work was supported in part by NSF I/UCRC grant 1747751.

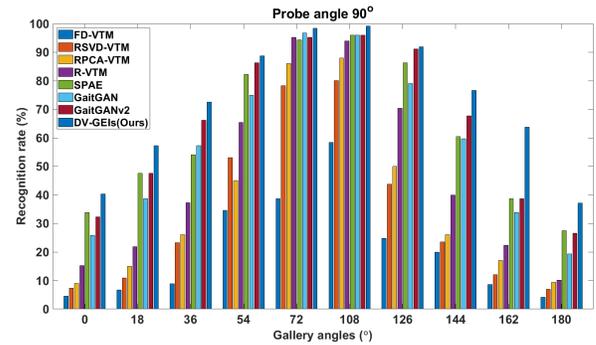
References

[1] W. An, R. Liao, S. Yu, Y. Huang, and P. C. Yuen. Improving gait recognition with 3d pose estimation. In *the 13th Chinese Conference on Biometric Recognition*, pages 137–147, 2018.

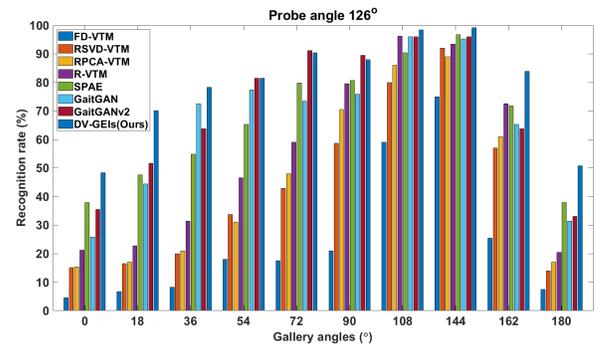
[2] W. An, S. Yu, Y. Makihara, X. Wu, C. Xu, Y. Yu, R. Liao, and Y. Yagi. Performance evaluation of model-based gait on multi-view very large population database with pose sequences. *IEEE Transactions on Biometrics, Behavior, and Identity Science*, 2020.



(a)



(b)



(c)

Figure 7. Comparisons with view transformation model methods at probe angles (a)54°, (b)90° and (c)126°. The gallery angles are the rest 10 angles except the corresponding probe angle. Our result outperforms those VTM methods.

[3] H. Chao, Y. He, J. Zhang, and J. Feng. Gaitset: Regarding gait as a set for cross-view gait recognition. In *Proceedings of the AAAI Conference on Artificial Intelligence*, pages 8126–8133, 2019.

[4] J. Chen, J. Chen, H. Chao, and M. Yang. Image blind denoising with generative adversarial network based noise modeling. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3155–3164, 2018.

[5] I. J. Goodfellow, J. Pougetabadié, M. Mirza, B. Xu, D. Wardefarley, S. Ozair, A. Courville, Y. Bengio,

- Z. Ghahramani, and M. Welling. Generative adversarial nets. *Advances in Neural Information Processing Systems*, 3:2672–2680, 2014.
- [6] J. Han and B. Bhanu. Individual recognition using gait energy image. *IEEE Transactions on Pattern Analysis & Machine Intelligence*, 28(2):316–22, 2006.
- [7] X. Hou, L. Shen, K. Sun, and G. Qiu. Deep feature consistent variational autoencoder. In *2017 IEEE Winter Conference on Applications of Computer Vision*, pages 1133–1141, 2017.
- [8] P. Isola, J.-Y. Zhu, T. Zhou, and A. A. Efros. Image-to-image translation with conditional adversarial networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1125–1134, 2017.
- [9] H. Iwama, M. Okumura, Y. Makihara, and Y. Yagi. The ou-isir gait database comprising the large population dataset and performance evaluation of gait recognition. *IEEE Transactions on Information Forensics and Security*, 7(5):1511–1521, 2012.
- [10] W. Kusakunniran, Q. Wu, H. Li, and J. Zhang. Multiple views gait recognition using view transformation model based on optimized gait energy image. In *ICCV Workshops*, pages 1058–1064, 2009.
- [11] W. Kusakunniran, Q. Wu, H. Li, and J. Zhang. Multiple views gait recognition using view transformation model based on optimized gait energy image. In *IEEE International Conference on Computer Vision Workshops*, pages 1058–1064, 2010.
- [12] W. Kusakunniran, Q. Wu, J. Zhang, and H. Li. Gait recognition under various viewing angles based on correlated motion regression. *IEEE TCSVT*, 22(6):966–980, 2012.
- [13] R. Liao, C. Cao, E. B. Garcia, S. Yu, and Y. Huang. Pose-based temporal-spatial network (ptsn) for gait recognition with carrying and clothing variations. In *the 12th Chinese Conference on Biometric Recognition*, pages 474–483, 2017.
- [14] R. Liao, S. Yu, W. An, and Y. Huang. A model-based gait recognition method with body pose and human prior knowledge. *Pattern Recognition*, 98:107069, 2020.
- [15] Y. Makihara, R. Sagawa, Y. Mukaigawa, T. Echigo, and Y. Yagi. Gait recognition using a view transformation model in the frequency domain. In *Proceedings of the European Conference on Computer Vision*, pages 151–163, 2006.
- [16] Y. Makihara, R. Sagawa, Y. Mukaigawa, T. Echigo, and Y. Yagi. Gait recognition using a view transformation model in the frequency domain. In *ECCV*, pages 151–163, 2006.
- [17] X. Qian, Y. Fu, T. Xiang, W. Wang, J. Qiu, Y. Wu, Y.-G. Jiang, and X. Xue. Pose-normalized image generation for person re-identification. In *Proceedings of the European Conference on Computer Vision*, pages 650–667, 2018.
- [18] O. Ronneberger, P. Fischer, and T. Brox. U-net: Convolutional networks for biomedical image segmentation. In *International Conference on Medical image computing and computer-assisted intervention*, pages 234–241. Springer, 2015.
- [19] N. Takemura, Y. Makihara, D. Muramatsu, T. Echigo, and Y. Yagi. Multi-view large population gait dataset and its performance evaluation for cross-view gait recognition. *IPSP Transactions on Computer Vision and Applications*, 10(1):4, 2018.
- [20] Y. Wen, K. Zhang, Z. Li, and Y. Qiao. A discriminative feature learning approach for deep face recognition. In *European conference on computer vision*, pages 499–515, 2016.
- [21] Z. Wu, Y. Huang, L. Wang, X. Wang, and T. Tan. A comprehensive study on cross-view gait based human identification with deep cnns. *IEEE Transactions on Pattern Analysis & Machine Intelligence*, 39(2):209–226, 2017.
- [22] S. Yu, H. Chen, E. B. G. Reyes, and N. Poh. GaitGAN: Invariant gait feature extraction using generative adversarial networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*, pages 30–37, 2017.
- [23] S. Yu, R. Liao, W. An, H. Chen, E. B. G. Reyes, Y. Huang, and N. Poh. GaitGANv2: Invariant gait feature extraction using generative adversarial networks. *Pattern recognition*, 87:179–189, 2019.
- [24] S. Yu, D. Tan, and T. Tan. A framework for evaluating the effect of view angle, clothing and carrying condition on gait recognition. In *the 18th International Conference on Pattern Recognition*, pages 441–444, 2006.
- [25] S. Yu, Q. Wang, L. Shen, and Y. Huang. View invariant gait recognition using only one uniform model. In *International Conference on Pattern Recognition*, pages 889–894, 2017.
- [26] S. Zheng, J. Zhang, K. Huang, R. He, and T. Tan. Robust view transformation model for gait recognition. In *IEEE International Conference on Image Processing*, pages 2073–2076, 2011.
- [27] S. Zheng, J. Zhang, K. Huang, R. He, and T. Tan. Robust view transformation model for gait recognition. In *ICIP*, pages 2073–2076, 2011.