



OPEN ACCESS

EDITED BY
Qiong Wu,
Jiangnan University, China

REVIEWED BY
Xiaobo Wang,
Jiangnan University, China
Hai Li,
II Aviation Flight University of China,
China

*CORRESPONDENCE
Jing Fan,
fanjing9476@163.com

SPECIALTY SECTION
This article was submitted to Aerial and
Space Networks,
a section of the journal
Frontiers in Space Technologies

RECEIVED 08 August 2022
ACCEPTED 12 September 2022
PUBLISHED 27 September 2022

CITATION
Tang Y, Fan J and Qu J (2022), High-
quality facial-expression image
generation for UAV
pedestrian detection.
Front. Space Technol. 3:1014183.
doi: 10.3389/frspt.2022.1014183

COPYRIGHT
© 2022 Tang, Fan and Qu. This is an
open-access article distributed under
the terms of the [Creative Commons
Attribution License \(CC BY\)](https://creativecommons.org/licenses/by/4.0/). The use,
distribution or reproduction in other
forums is permitted, provided the
original author(s) and the copyright
owner(s) are credited and that the
original publication in this journal is
cited, in accordance with accepted
academic practice. No use, distribution
or reproduction is permitted which does
not comply with these terms.

High-quality facial-expression image generation for UAV pedestrian detection

Yumin Tang^{1,2}, Jing Fan^{1,2*} and Jinshuai Qu^{1,2}

¹School of Electrical Information Engineering, Yunnan Minzu University, Kunming, China, ²University Key Laboratory of Information and Communication on Security Backup and Recovery in Yunnan Province, Yunnan Minzu University, Kunming, China

For UAV pedestrian detection in the wild with perturbed parameters, such as lighting, distance, poor pixel and uneven distribution, traditional methods of image generation cannot accurately generate facial-expression images for UAV pedestrian detection. In this study, we propose an improved PR-SGAN (perceptual-remix-star generative adversarial network) method, which combines the improved interpolation method, perceptual loss function, and StarGAN to achieve high-quality facial-expression image generation. Experimental results show that the proposed method for discriminator-parameter update improves the generated facial-expression images in terms of image-generation evaluation indexes (5.80 dB in PSNR and 24% in SSIM); the generated images for generator-parameter update have high robustness against color. Compared to the traditional StarGAN method, the generated images are significantly improved in high frequency details and textures.

KEYWORDS

facial expression recognition, generative adversarial network, StarGAN, high quality, high robustness

1 Introduction

With the development of computer vision and the field of UAVs (unmanned aerial vehicles), target detection based on UAV has gained more and more attention in the military and civilian fields. With the convenience and speed of traffic (Wu et al., 2022a), the escape range and speed of criminals are also expanding, so real-time monitoring becomes more necessary and important (Srivastava et al., 2022). Compared with traditional camera monitoring, UAV have better flexibility and operability, making it possible to monitor any range of anomalies and provide real-time feedback.

Based on this, the research on pedestrian detection based on UAV become more and more important. However, due to the different viewing angles of UAV, UAV pedestrian detection has many characteristics, such as complex backgrounds, small-scale targets, uneven distribution, etc. (Wang et al., 2022). Pedestrian detection based on UAV is mainly for the accurate identification of the target, while the precise positioning of pedestrians is mainly for face recognition, and the accuracy of the facial recognition rate mainly depends on the amount of data and the corresponding improvement of image

quality. The generation of high-quality facial images in pedestrian detection is of great help to improve the accuracy of UAV pedestrian detection.

The generation of facial expression images is a more refined image generation method for face image generation, which can make the generated images more realistic, and can make the generated images more appropriate to the original face images. After the UAV pedestrian detection completes the contour positioning of the person, the accuracy of the positioning and detection of the face is greatly improved. However, most of the recent researches on the generation of facial expression images blindly aim to improve the quality of image generation in terms of image indicators, and seldom pay attention to the improvement of the perceptual effect of facial expression image generation. As a result, the quality of the generated facial expression image has a good effect on the index, but the actual perception effect is very poor. Consequently, the facial recognition accuracy of the generated face image cannot be well improved.

For this paper, the main contributions are as follows:

- 1) A method that achieves simultaneous improvement in image generation quality and expansion of image quantity is proposed and applied to the face expression image generation method for UAV pedestrian detection with high robustness and reliability.

- 2) Multiple experiments are conducted on a facial expression dataset in the wild to verify the achievability of the method in UAV pedestrian detection.

The remainder of this paper consists of the following contents. Section 2 summarizes related work. Section 3 gives a detailed description of the content of the method proposed in this paper. Section 4 explains the system model of the proposed method. Section 5 provides a detailed analysis of the experiments and experimental results. Section 6 concludes the paper.

2 Related works

2.1 Image generation for UAV pedestrian detection

With the rapid development of image processing technology, corresponding fields such as computer vision have also made great progress. The continuous progress of deep learning methods also enables it to develop more complex algorithms and apply them to UAV tasks, such as infrastructure monitoring (Banić et al., 2019; Peng et al., 2021), crop analysis (Banerjee et al., 2021; Donmez et al., 2021), crowd counting method in UAV pedestrian detection (Ptak and Pieczynski, 2022) (calculating people collected by UAV Low Altitude Aerial Photography),

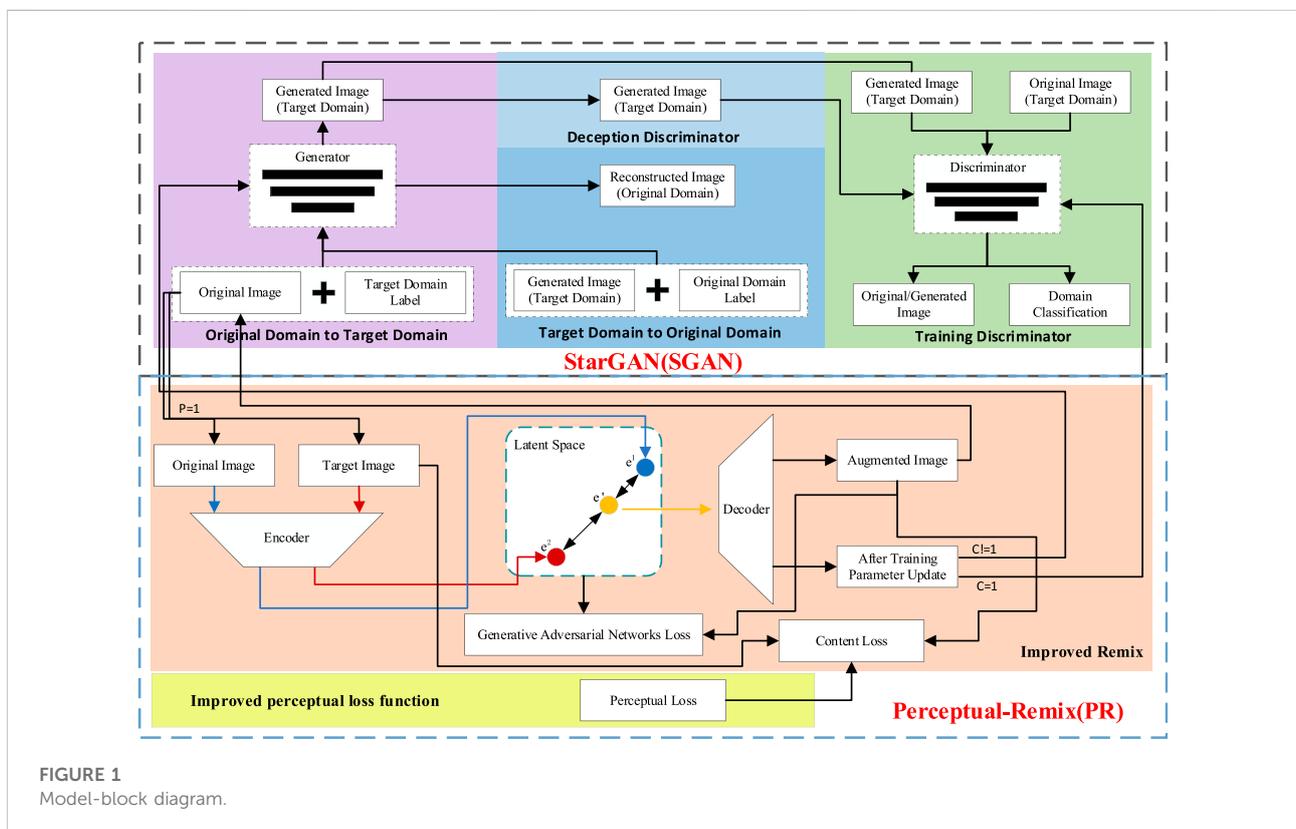


FIGURE 1 Model-block diagram.

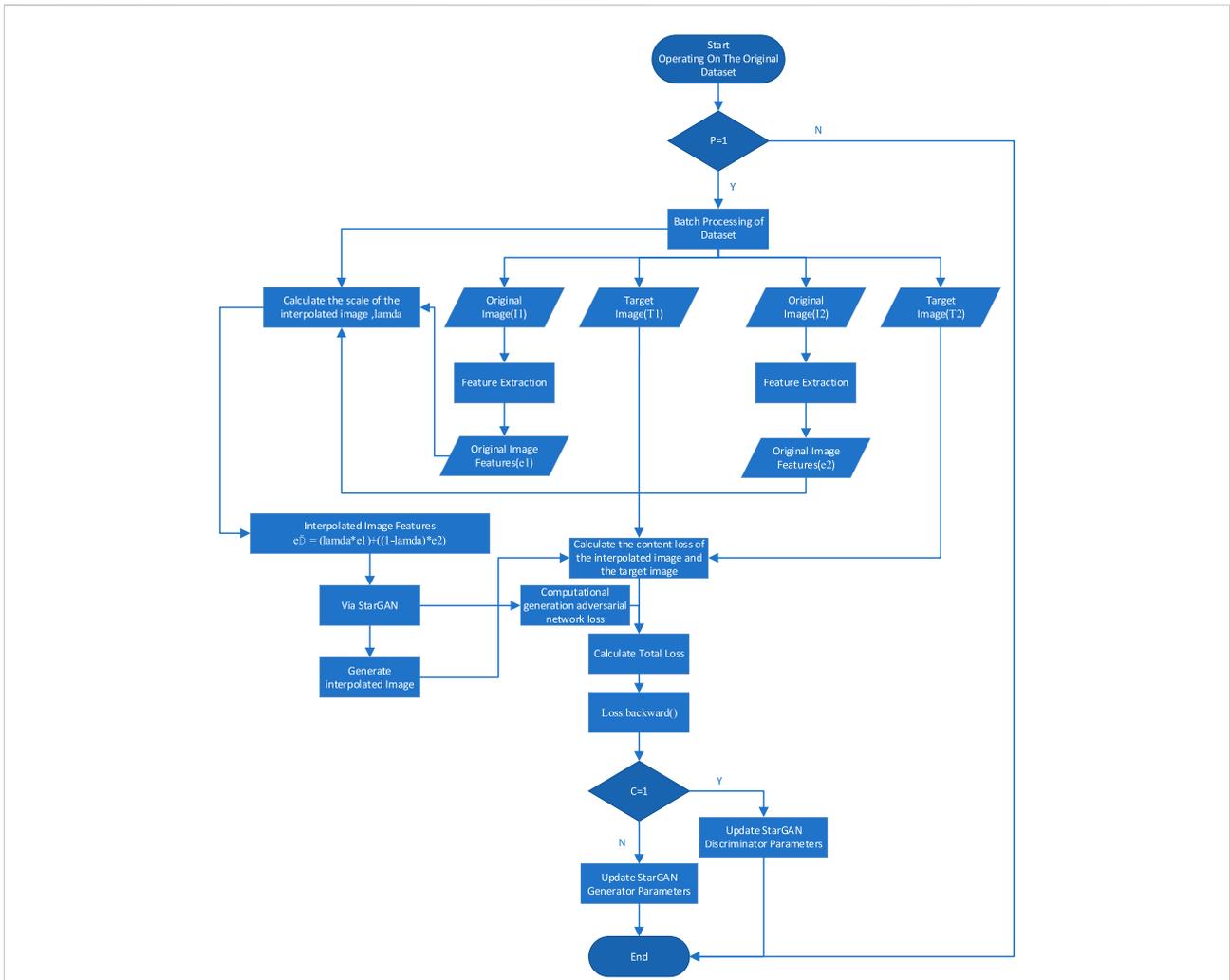


FIGURE 2
Flowchart of the improved Remix algorithm.

autonomous driving (Wu and Zheng, 2015; Wu and Zheng, 2016; Wu et al., 2022b; Zhu et al., 2022), and so on. Ptak and Pieczynski (2022) used the tools provided by the game engine to generate a crowd count dataset with UAV view characteristics, which can realize the image flow of the simulator and the active control of the UAV. Park et al. (2022) proposed a (range-doppler) RD data augmentation method based on conditionally generated adversarial network to solve the problem of lack of UAV related task datasets; according to the collected UAV RD map, a synthetic UAV RD map was generated, and a new UAV RD map dataset was generated; The experimental results verified the effectiveness of the image expansion method.

At present, although many studies have made great breakthroughs in target detection, high-resolution UAV target detection is still a very challenging research due to the special characteristics of the images collected by UAV. There are three

main reasons for this: 1) The objects in the data collected by UAV are usually small; 2) The size of the image collection is too large, and it needs to be compressed when placed in the detection technology, so that the effect of the image is not as good as the original image; 3) The data collected by UAV is unevenly distributed. This paper will solve the problem of poor quality of facial expression image generation in UAV pedestrian detection to achieve image quality improvement.

2.2 Loss function

Ledig et al. (2017) posited that MSE (Mean Square Error) makes the generated images considerably smooth, resulting in poor perceptual quality level of the generated images.

In the facial-expression recognition task, the poor perceptual quality of the generated images has a detrimental impact on the

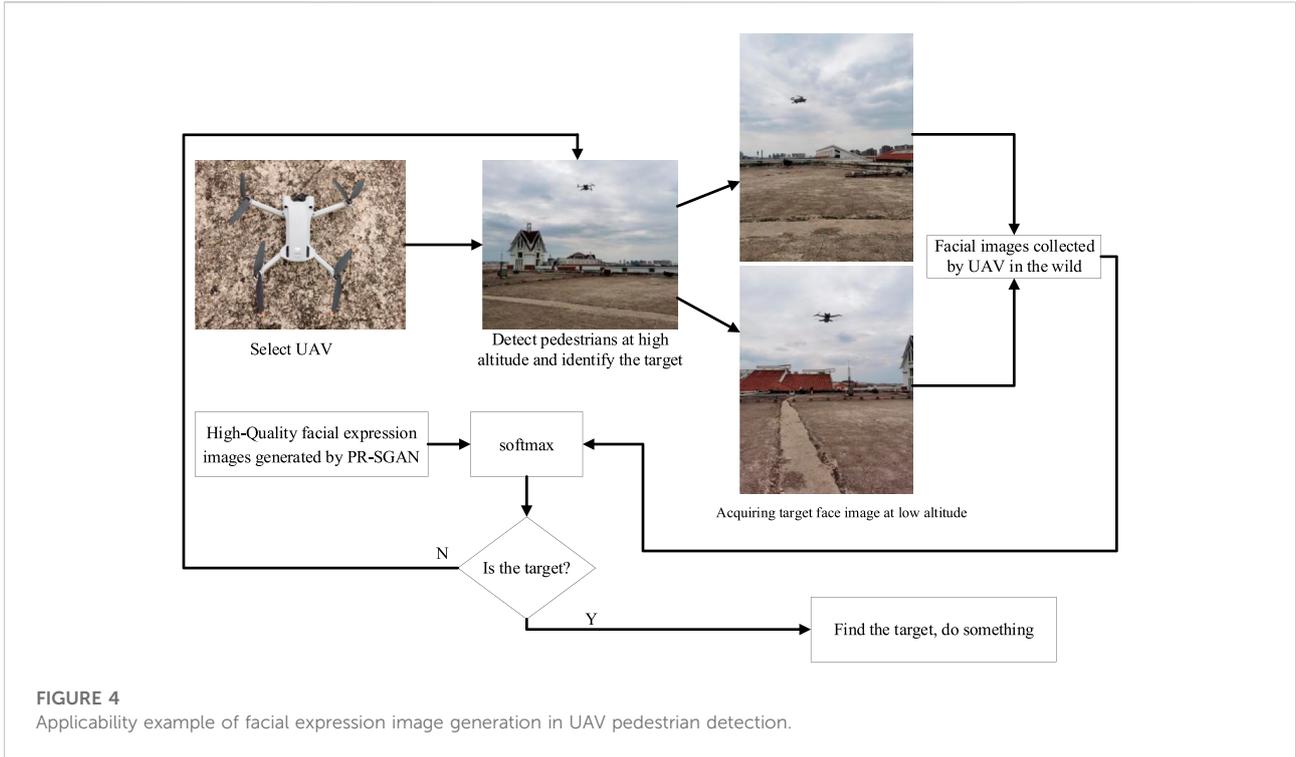


FIGURE 4
Applicability example of facial expression image generation in UAV pedestrian detection.

acquired images. This is because MSE is averaged over pixels and can make the image extremely smooth. Conversely, the goal of high-quality image generation is the visual fidelity of the generated fake images. To improve the effect of facial-expression images, the MSE loss function in (Johnson et al., 2016) is improved to a perceptual loss function applicable to the improved network in this study. The loss function formulas after the improvement are expressed in Eqs 6, 7. The perceptual loss model is a VGG19 network using the number of layers with activation values of the conv1-5 modules. The weights are assigned for feature matching of each module, where φ is the loss network, $C_{i,j,k}H_{i,j,k}W_{i,j,k}$ is the dimensional features in the VGG network, $I_{x,y,z}$ is the original image, $T_{x,y,z}$ is the target image, and R is the improved Remix algorithm. The improved loss function achieves the conversion of the computed space from the pixel space to the feature space, and improves the image-generation effect.

$$L_{P1} = \frac{1}{C_{i,j,k}H_{i,j,k}W_{i,j,k}} \sum_{x=1}^{C_{i,j,k}} \sum_{y=1}^{H_{i,j,k}} \sum_{z=1}^{W_{i,j,k}} (\varphi_{i,j,k}(T_{1x,y,z}) - \varphi_{i,j,k}(R(I_{1x,y,z}, I_{2x,y,z})))^2 \quad (6)$$

$$L_{P2} = \frac{1}{C_{i,j,k}H_{i,j,k}W_{i,j,k}} \sum_{x=1}^{C_{i,j,k}} \sum_{y=1}^{H_{i,j,k}} \sum_{z=1}^{W_{i,j,k}} (\varphi_{i,j,k}(T_{2x,y,z}) - \varphi_{i,j,k}(R(I_{1x,y,z}, I_{2x,y,z})))^2 \quad (7)$$

3.4 Improved generator network structure

Because of the Remix interpolation method, when updating the parameters of the generator, the image becomes overly

sensitive to the color and the generated facial-expression image is not robust to color. We propose a method for improving the network structure of the generator by replacing the ReLU activation function in the original generator network with the more stable PReLU (He et al., 2015) activation function. Figure 3 shows the improved network structure.

4 System model

As shown in Figure 4, it is an application example of facial expression image generation in UAV pedestrian detection. The specific steps are: 1) Determine the pedestrian object to be detected; 2) Determine the UAV to collect the target object; 3) The UAV flies to a high altitude and locates the pedestrian whose overall outline conforms; 4) The UAV flies low at the positioning position and collects the facial expression images of pedestrians; 5) The collected facial expression images and the images generated by using the PR-SGAN algorithm are used for softmax binary classification; 6) If the object correspondence is accurate, the accurate pedestrian detection of the UAV is completed; if the correspondence is not accurate, go back to the step 3 for pedestrian detection by UAV. The research on this aspect can greatly promote the progress of UAV pedestrian detection and further promote its application in the Internet of things (Wu et al., 2019).



5 Experiments and results

5.1 Experimental dataset

In this study, the images are first cropped to 256×256 and scaled to 128×128 to ensure that the borders of the facial expression image images in the RAF-DB (Li et al., 2017) dataset are not distorted.

5.2 Experimental environment and parameter settings

This experiment was conducted using the Centos operating system. We used a mini-batch with a value of 16 to improve the

efficiency of model training. We set the learning rates of the generator and discriminator to 0.0001, and the number of learning-rate decay iterations to 10,000. We set the number of iterations for the training process and the discriminator to 200 and 520 K, respectively.

Additionally, we set the properties in the interpolation algorithm of the Remix method to one, and the momentum to 0.9. Compared to the original code (Cao and Hou Yang He, 2021), the improved code is more flexible. By setting the value of the P, we could decide whether to use the Remix interpolation expansion. Simultaneously, we used C to realize the selections for updating the generator and discriminator parameters. When C is one, the discriminator is updated; otherwise, it is not updated.

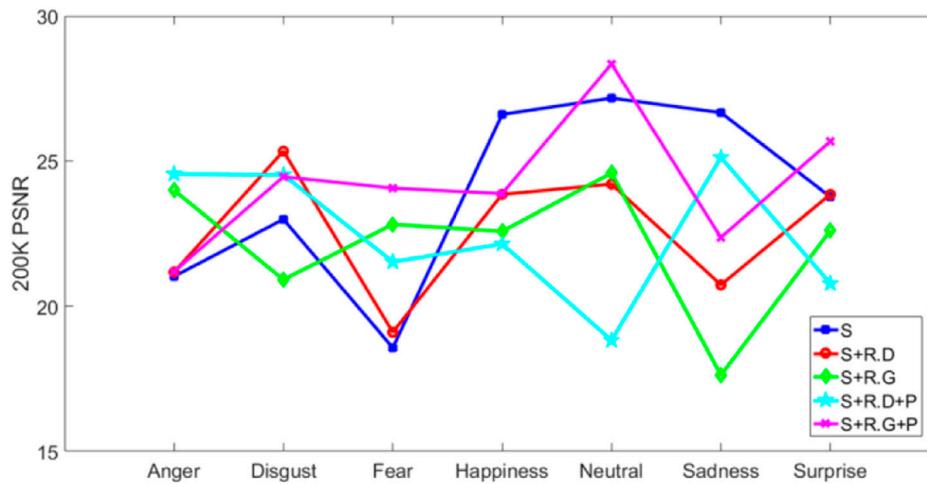


FIGURE 6
Number of iterations is 200 K the SSIM value of the generated image.

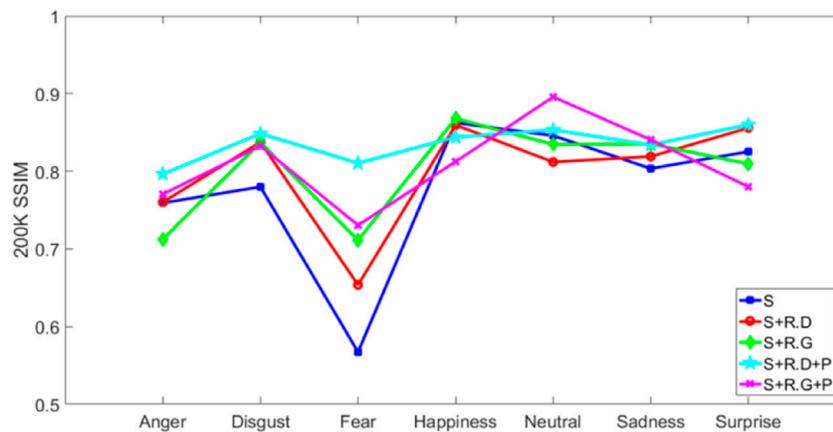


FIGURE 7
Number of iterations is 200 K the PSNR value of the generated image.

5.3 Experimental effect comparison image

In this comparison test, seven facial images were randomly selected as input to generate the corresponding seven expression images. Because expression classification focuses on features, such as eyes and lips, we compare the effects generated by the original code with the images generated by the improved method on these features.

Figure 5 shows the comparison of the generated facial-expression images with 200 K iterations. Where S represents StarGAN and R.D and R.G are the update steps in the Remix code. R.D opts to update the discriminator and R.G selects the generator; P refers to the perceptual loss function. When

updating the discriminator based on the original Remix method, facial noise and perturbation increase, ripples and other colors appear in the image, and partial distortion occurs in the lips and eyes, which is not promising for high-quality image generation; simultaneously, ghosting appears when updating the generator. For the above problems, we add the perceptual loss function. From the comparison graph, the generated image of S + R.D + P is more vivid and complete in the part of lips and eyes when the number of iterations is 200 K and the expression changes obviously and the artifact noise is eliminated. The effect of S + R.G + P-generated lips and eyes is a slightly worse compared to that of S + R.D + P; however, it is better than that of S + R.G and S.

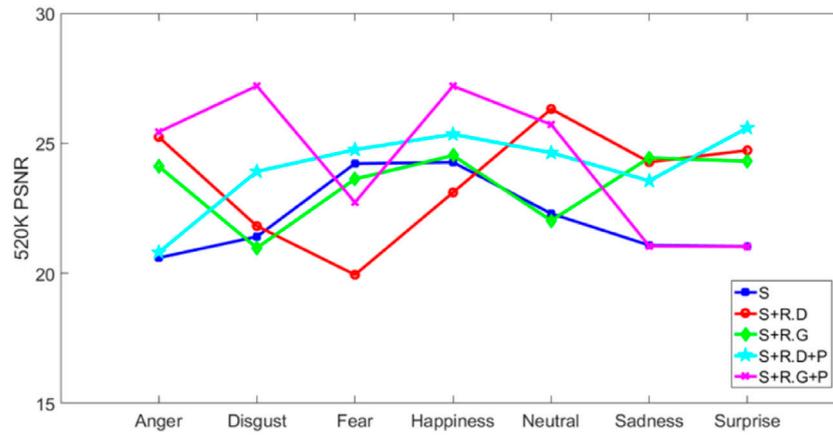


FIGURE 8
Number of iterations is 520 K the PSNR value of the generated image.

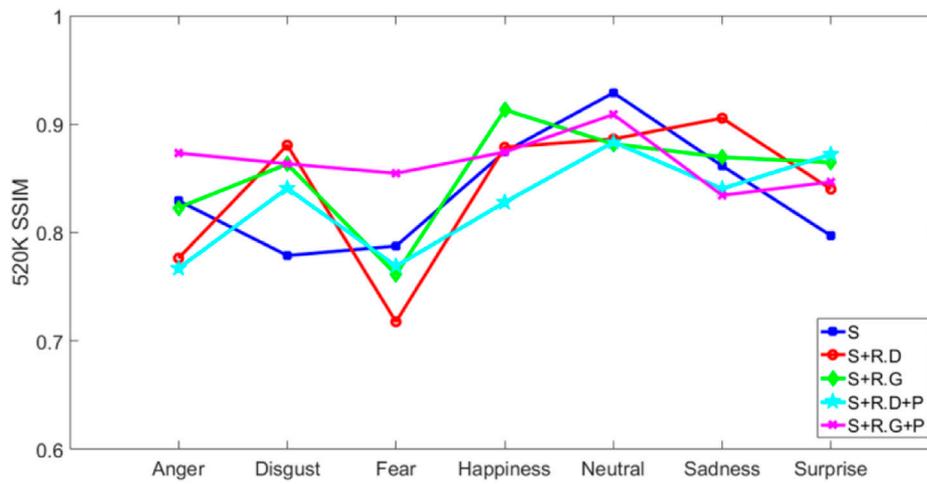


FIGURE 9
Number of iterations is 520 K the SSIM value of the generated image.

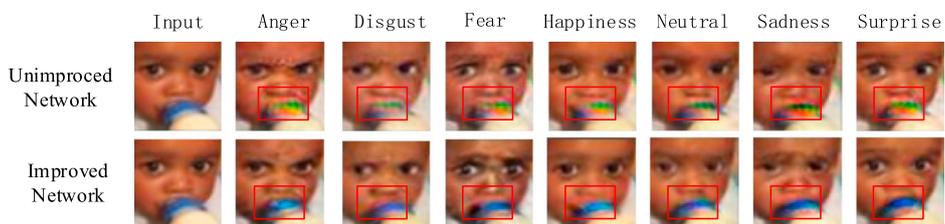


FIGURE 10
Facial-expression image generated by improved and unimproved networks.

Figure 5 shows that when the number of iterations is 520 K the effect of eye and mouth generation is affected by the discriminator-parameter update in the Remix-based method; however, the ripple and color interference disappears, such that the image-generation effect is not promising. The update for the generator is affected by a similar degree of interference and influence. By adding the perceptual loss function for the above problem, the image generated by S + R.D + P becomes more vivid and the effect of the mouth and eyes on the corresponding expressions disappears. The affected effect is reduced, and the generated image obtained is more realistic and reliable. The repair effect of S + R.G + P on the features is better than that of S + R.G and S.

Generally, as the number of iterations continues to improve, the generated image effect also improves, and the generated image effect of S + R.D + P is the most excellent. The improved method based on the generator is not sufficiently robust for the color aspect, and the subsequent improvement of this study addresses this problem and conducts a comparison test.

5.4 Indicators results

The PSNR and SSIM values for this experiment are obtained by calculating the PSNR and SSIM values from a set of generated images with the original images in matlab 2016a environment. These experimental data graphs respectively illustrate the new seven corresponding expression images generated by using the corresponding method for the original image under different iteration times, and calculate the image quality discrimination index between the generated image and the original image, so as to better compare the differences in indexes of the images generated by different methods.

As shown in Figures 6–9, most of the face images generated by the our proposed method are of higher quality compared to those generated by the original StarGAN in the two metrics of PSNR and SSIM. The comparison on the indices shows the feasibility and realism of the improved method proposed in this study.

Specifically, as shown in Figure 6, under 200 K iterations, the number of values greater than S and the number less than S of the improved method in this paper are almost the same; while Figure 8 shows that at 520 K iteration, the number greater than S is one more than the number greater than S. In addition, according to Figure 6 and Figure 8, it can be seen that only the facial expression images with two expressions in the s method are in the leading position in terms of indicators, while the remaining expressions are in the leading position of the method proposed in this paper, and have achieved a good effect in improving indicators in the expression of “disgust,” which better reflects the excellence of the method proposed in this paper. In short, the method proposed in this paper is effective and reliable in improving SSIM value.

In terms of PSNR value, as shown in Figure 7 and Figure 9 the image PSNR value obtained by the method proposed in this

paper at 200 and 520 K iterations is significantly more than S. In addition, the image generation effect of the proposed method in all expressions are better than the S method. It can be seen that the method proposed in this paper also has strong reliability and effectiveness in improving the PSNR index of images.

5.5 Experimental results of the improved generator network structure

Because poor color robustness occurs when updating parameters to the generator, the network of the generator is improved for the case where C is not taken as one. Experiments prove the effectiveness of the improvement of the generator-network structure. Figure 10 shows the comparison of the facial-expression images generated by the unimproved and improved networks with 520 K iterations based on the updated parameters of the generator. Evidently, after the improvement of the generator network, the generated facial-expression images are more robust in color and smoother in texture, and the quality of the generated images is improved.

6 Conclusion

In this paper, we propose an improved PR-SGAN method. It can independently generate high-quality face expression images of a person with different expressions by learning the relationship between different domains in the dataset. The effectiveness of our method is verified on the comparative experiments of discriminative indexes and image generation effects. It is praiseworthy that the PR-SGAN method is a good solution to the problem of over-fitting and poor quality of image generation for UAV pedestrian detection in the wild. In the future, more in-depth research will be conducted to generate higher quality facial-expression images to improve the accuracy and robustness for UAV pedestrian detection in the wild.

Data availability statement

The original contributions presented in the study are included in the article/supplementary material, further inquiries can be directed to the corresponding author.

Author contributions

YT completed the corresponding algorithm implementation and paper writing, JF provided guidance on the algorithm of the paper and the writing of the paper, and JQ provided corresponding guidance on the writing of the paper.

Funding

This research was funded by the National Natural Science Foundation of China (Grant No. 61540063), Application Basic Research Fund of Yunnan Province (Grant No. 2018FD055), MOE (Ministry of Education in China) Project of Humanities and Social Sciences (Grant No. 20YJCZH129) and Yunnan Provincial Department of Education Science Research Fund Project (Grant No. 2020J0655).

Acknowledgments

We would like to thank Editage (www.editage.cn) for the English language editing.

References

- Banerjee, B. P., Sharma, V., Spangenberg, G., and Kant, S. (2021). Machine learning regression analysis for estimation of crop emergence using multispectral UAV imagery. *Remote Sens.* 1315, 2918. doi:10.3390/rs13152918
- Banić, M., Miltenović, A., Pavlović, M., and Ćirić, I. (2019). Intelligent machine vision based railway infrastructure inspection and monitoring using UAV. *Facta Univ. Ser. Mech. Eng.* 173, 357–364.
- Cao, J., and Hou Yang He, L. M. H. R. (2021). “ReMix: Towards image-to-image translation with limited data,” in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. doi:10.1109/cvpr46437.2021.01477
- Choi, Y., Choi, M., Kim, M., and Ha, J. W. (2018). “Stargan: Unified generative adversarial networks for multi-domain image-to-image translation,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*. doi:10.1109/cvpr.2018.00916
- Donmez, C., Villi, O., Berberoglu, S., and Cilek, A. (2021). Computer vision-based citrus tree detection in a cultivated environment using UAV imagery. *Comput. Electron. Agric.* 187, 106273. doi:10.1016/j.compag.2021.106273
- He, K., Zhang, X., Ren, S., and Sun, J. (2015). “Delving deep into rectifiers: Surpassing human-level performance on imagenet classification,” in *Proceedings of the IEEE international conference on computer vision*. doi:10.1109/iccv.2015.123
- Johnson, J., Alexandre, A., and Fei-Fei, L. (2016). “Perceptual losses for real-time style transfer and super-resolution,” in *European conference on computer vision* (Cham: Springer). doi:10.1007/978-3-319-46475-6_43
- Ledig, C., Theis, L., Huszar, F., Caballero, J., Cunningham, A., Acosta, A., et al. (2017). “Photo-realistic single image super-resolution using a generative adversarial network,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*. doi:10.1109/cvpr.2017.19
- Li, S., Deng, W., and Du, J. P. (2017). “Reliable crowdsourcing and deep locality-preserving learning for expression recognition in the wild,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*. doi:10.1109/cvpr.2017.277
- Lucas, A., Lopez-Tapia, S., Molina, R., and Katsaggelos, A. K. (2019). Generative adversarial networks and perceptual losses for video super-resolution. *IEEE Trans. Image Process.* 28 (7), 3312–3327. doi:10.1109/tip.2019.2895768
- Park, S., Lee, S., and Kwak, N. (2022). “Range-Doppler map augmentation by generative adversarial network for deep UAV classification,” in *2022 IEEE radar conference (RadarConf22)* (IEEE). doi:10.1109/radarconf2248738.2022.9764177

Conflict of interest

The authors declare that the research was conducted in the absence of any commercial or financial relationships that could be construed as a potential conflict of interest.

Publisher’s note

All claims expressed in this article are solely those of the authors and do not necessarily represent those of their affiliated organizations, or those of the publisher, the editors and the reviewers. Any product that may be evaluated in this article, or claim that may be made by its manufacturer, is not guaranteed or endorsed by the publisher.

- Peng, X., Zhong, X., Zhao, C., Chen, A., and Zhang, T. (2021). A UAV-based machine vision method for bridge crack recognition and width quantification through hybrid feature learning. *Constr. Build. Mater.* 299, 123896.
- Ptak, B., and Pieczynski, D. (2022). *CountingSim: Synthetic way to generate a dataset for the UAV-view crowd counting task*.
- Srivastava, A., Badal, T., Saxena, P., Vidyarthi, A., and Singh, R. (2022). UAV surveillance for violence detection and individual identification. *Autom. Softw. Eng.* 291, 1–28. doi:10.1007/s10515-022-00323-3
- Wang, A., Fang, Z., Gao, Y., Jiang, X., and Ma, S. (2018). Depth estimation of video sequences with perceptual losses. *IEEE Access* 6, 30536–30546. doi:10.1109/access.2018.2846546
- Wang, C., Luo, D., Liu, Y., Xu, B., and Zhou, Y. (2022). Near-surface pedestrian detection method based on deep learning for UAVs in low illumination environments. *Opt. Eng.* 612, 023103. doi:10.1117/1.oe.61.2.023103
- Wu, Q., Fan, Y. Q., and Fan, Q. (2022). “Time-dependent performance modeling for platooning communications at intersection,” in *IEEE Internet things journal*, 1. doi:10.1109/JIOT.2022.3161028
- Wu, Q., Liu, H., Zhang, C., Fan, Q., Li, Z., and Wang, K. (2019). Trajectory protection schemes based on a gravity mobility model in IoT. *Electronics* 8, 148. doi:10.3390/electronics8020148
- Wu, Q., Wan, Z., Fan, Q., Fan, P., and Wang, J. (2022). Velocity-adaptive access scheme for MEC-assisted platooning networks: Access fairness via data freshness. *IEEE Internet Things J.* 9 (6), 4229–4244. doi:10.1109/jiot.2021.3103325
- Wu, Q., and Zheng, J. (2015). “Performance modeling and analysis of the ADHOC MAC protocol for VANETs,” in *2015 IEEE international conference on communications (ICC)* (IEEE). doi:10.1109/icc.2015.7248891
- Wu, Q., and Zheng, J. (2016). Performance modeling and analysis of the ADHOC MAC protocol for vehicular networks. *Wirel. Netw.* 223, 799–812. doi:10.1007/s11276-015-1000-6
- Zhang, X., Ren, N., and Chen, Q. (2018). “Single image reflection separation with perceptual losses,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*. doi:10.1109/cvpr.2018.00503
- Zhu, H., Wu, Q., Wu, X.-J., Fan, Q., Fan, P., and Wang, J. (2022). Decentralized power allocation for MIMO-NOMA vehicular edge computing based on deep reinforcement learning. *IEEE Internet Things J.* 9 (4), 12770–12782. doi:10.1109/jiot.2021.3138434