

Generative Active Learning with Variational Autoencoder for Radiology Data Generation in Veterinary Medicine

In-Gyu Lee

*Dept. Computer Science
Chungbuk National University
Cheongju, Republic of Korea
ingyu.lee@chungbuk.ac.kr*

Jun-Young Oh

*Dept. Computer Science
Chungbuk National University
Cheongju, Republic of Korea
jy.oh@chungbuk.ac.kr*

Hee-Jung Yu

*Dept. Veterinary Medical Imaging
Konkuk University
Seoul, Republic of Korea
hazel13@konkuk.ac.kr*

Jae-Hwan Kim

*Dept. Veterinary Medical Imaging
Konkuk University
Seoul, Republic of Korea
jaehwan@konkuk.ac.kr*

Ki-Dong Eom

*Dept. Veterinary Medical Imaging
Konkuk University
Seoul, Republic of Korea
eomkd@konkuk.ac.kr*

Ji-Hoon Jeong*

*Dept. Computer Science
Chungbuk National University
Cheongju, Republic of Korea
jh.jeong@chungbuk.ac.kr*

Abstract—Recently, with increasing interest in pet healthcare, the demand for computer-aided diagnosis (CAD) systems in veterinary medicine has increased. The development of veterinary CAD has stagnated due to a lack of sufficient radiology data. To overcome the challenge, we propose a generative active learning framework based on a variational autoencoder. This approach aims to alleviate the scarcity of reliable data for CAD systems in veterinary medicine. This study utilizes datasets comprising cardiomegaly radiograph data. After removing annotations and standardizing images, we employed a framework for data augmentation, which consists of a data generation phase and a query phase for filtering the generated data. The experimental results revealed that as the data generated through this framework was added to the training data of the generative model, the frechet inception distance consistently decreased from 84.14 to 50.75 on the radiograph. Subsequently, when the generated data were incorporated into the training of the classification model, the false positive of the confusion matrix also improved from 0.16 to 0.66 on the radiograph. The proposed framework has the potential to address the challenges of data scarcity in medical CAD, contributing to its advancement.

Index Terms—Artificial intelligence, generative model, active learning, variational autoencoder, data augmentation

I. INTRODUCTION

Pets have become important members of our lives, forming strong bonds and emotional connections with their owners. The healthcare of pets has become a subject of increased interest. With the continuous evolution of artificial intelligence (AI), there is a noticeable trend towards incorporating AI into computer-aided diagnosis (CAD) systems for pet healthcare. The effectiveness of AI models is heavily dependent on access to high-quality training data [1], [2]. However, acquiring a substantial amount of medical data for CAD has challenges due to the sensitive personal information in such data. Conse-

quently, there is a persistent effort to explore the application of generative models for the creation of medical data.

Among these efforts, numerous studies have leveraged generative adversarial networks (GAN) [3]. Yoon et al. [4] achieved a frechet inception distance (FID) of 42.19 for Sessile serrated lesion images using style-based GAN. Salvia et al. [5] proposed the use of GAN to generate synthetic hyperspectral images of epidermal lesions, addressing the challenge of limited large datasets. These academic efforts demonstrate the potential of GAN in generating medical images to improve research and diagnostics.

Zhu et al. [6] explored the concept of generative adversarial active learning (GAAL), employing a generative model in active learning to improve the performance of classification models. Although this approach involved queries to augment the training dataset for labeling, they emphasized the proposed framework, rather than focusing on criteria for queries.

In this study, we have creatively used active learning in the context of generative models, incorporating a unique criterion for data filtration through a variational autoencoder (VAE) [7]. This approach has significantly enhanced the robustness of the generative model's performance. We focused on addressing the scarcity of medical data for CAD, especially in the field of veterinary medicine, which may be helpful in medical AI advances. The introduced approach, termed Generative Active Learning with VAE, leverages query processes facilitated by a VAE to improve the performance of generative models in generating medical image data. This method stands out as a viable solution to address the persistent challenge of limited medical image data in CAD applications.

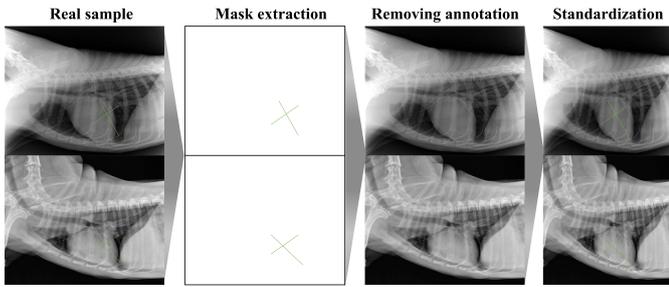


Fig. 1. The data preprocessing pipeline for training a generative model. If doctors draw annotations for diagnosis, the annotations are extracted to create masks. Then, image inpainting techniques are applied to remove the annotations. Subsequently, the resolution of the images is standardized.

II. GENERATIVE ACTIVE LEARNING WITH VAE

A. Dataset

In this study, we utilized the “Image Data for Diagnosis of Pet Diseases (thorax)” from AIHUB, a public dataset. We selected 100 images from cardiomegaly disease data. Data on cardiomegaly disease can visually confirm that the heart is enlarged [8]. We used these selected data as initial training data for the generative model.

Before training the model, a preprocessing step was performed to improve the quality of the data used for training. For radiographic images, annotations made by veterinarians for diagnosis were present. As the model could potentially learn from these annotations as features of the data, we performed a task to remove the annotations. Initially, the data in RGB format was transformed into the HSV color space. Subsequently, since the color of the annotations was in kinds of green, we extracted the green tones to create a binary mask. The mask obtained was then utilized in the image inpainting technique to restore the image. This method replaces pixels in the masked area using neighboring pixels.

Second, we standardized the resolution of both radiographic images. The raw data had diverse resolutions. Inconsistent image resolutions in the training dataset can lead to unstable learning due to variations in the size of the feature maps extracted by the neural network. To address this, we employ the center-cropping method, which uses center-based image cropping to ensure that essential organ information, such as the heart and kidney, is not lost. Radiographic images were resized to 256×256 pixels. The data preprocessing procedure is depicted in Fig. 1.

B. Proposed Framework

The framework is composed of the two phases. First, the data generating phase trains the generative model and generates data. Second, the query phase filters the generated data through the query strategy before incorporating them into the dataset for training the generative model. Algorithm 1 details the specific steps involved in this process with data preprocessing.

Algorithm 1 Overall process of proposed framework

Step 1: Data Preprocessing

Input: *Raw_dataset*

Output: *Preprocessed_dataset*

```

1 if Annotations exist in Raw_dataset then
2   Convert RGB images to HSV images
   Extract green tones to create masks
   Inpainted_dataset = Use masks for image inpainting
3 Preprocessed_dataset = Standardize resolution of Inpainted_dataset

```

Step 2: Data Generation

Input: *Preprocessed_dataset*

Output: *Augmented_dataset*

```

4 while Augmented_dataset size < 500 do
5   for epochs = 1 to 20 do
6     if epochs is even then
7       New_FID = Evaluate the generative model
8       if Saved_FID > New_FID then
9         Saved_FID = New_FID
10        Saved_weights = Save the generative model's weights
11      Generated_dataset = Generate 1000 data using Saved_weights
12      Calculate the cosine similarity of Generated_dataset
      Filtered_dataset = Select the top 10% Generated_dataset
      Augmented_dataset += Filtered_dataset

```

1) *Data generating phase:* The overall flow is depicted in Fig. 2. We refer to this entire process as a cycle and repeat the cycle 4 times until the train dataset has 500 data. The study utilized the projectedGAN model proposed by SaueI et al. [9] due to its state-of-the-art performance across various datasets during the experimentation. ProjectedGAN comprises a generator and a discriminator. The generator is trained to learn the data features to generate images that can effectively deceive the discriminator. The discriminator learns the data in a way that distinguishes between real data and generated data. The evaluation of the generation model was based on the FID [10], which measures the dissimilarity of the characteristics between the generated and actual images. A lower FID value indicates superior performance. The formula to calculate the FID is provided below. T represents actual images, and G represents generated images. Tr is defined as the sum of elements from the upper left to the lower right of the vector.

$$FID = \|\mu_T - \mu_G\|^2 - Tr((\Sigma_T + \Sigma_G - 2(\Sigma_T \Sigma_G)^{\frac{1}{2}})) \quad (1)$$

The GAN initiates training by using 100 selected actual data in the initial cycle. Training progresses through a total of 20 epochs per cycle, with a performance evaluation conducted every 2 epochs. Consequently, each cycle yields a total of 10 FID assessments. During evaluation, if the current FID value is lower than the previously recorded FID value, we save the model’s weights. We utilize these saved weights to generate 1,000 images per cycle.

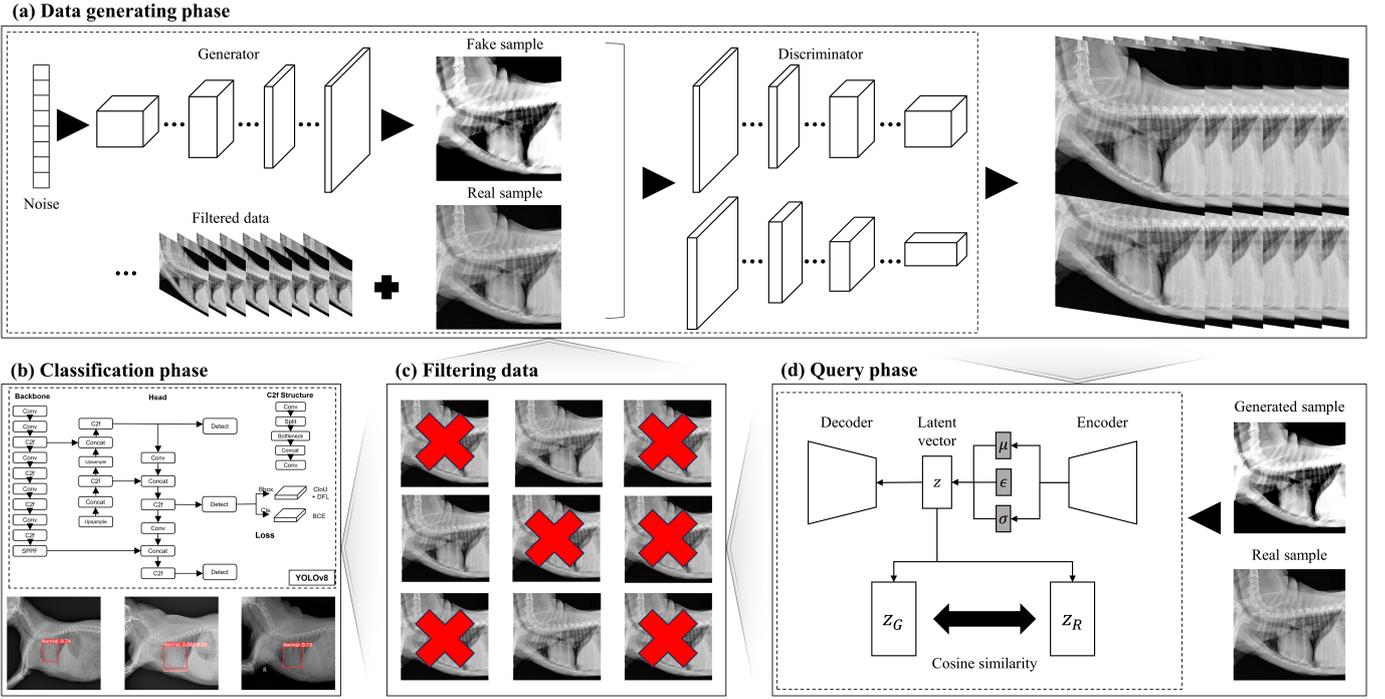


Fig. 2. Overall flow of the proposed framework. (a) The projectedGAN is trained with filtered image data and real image data to generate a new radiographic image. (d) VAE trained with 100 original data are used to filter generated images using a query strategy. (c) The top 10% cosine similarity of the data is added to the training dataset. (b) Finally, classification is performed using the object detection model after labeling to prove the usefulness of the data.

2) *Query phase*: The evaluation of image similarity in this study used a VAE, comprising an encoder and a decoder. In VAE, the decoder aims to regenerate the input in a form that is most similar to when a latent vector is given. The encoder, on the other hand, seeks to find the mean and standard deviation of the input and generates a latent vector with noise epsilon in Gaussian distribution. The study focused on utilizing the latent vector generated by the encoder. Training the VAE involved using 100 selected actual data for a total of 25 epochs. Training the VAE involved using 100 selected actual data for a total of 25 epochs.

To assess image similarity, cosine similarity was employed [11]. Unlike distance measures such as the Euclidean distance, which evaluate vectors based on their magnitudes, cosine similarity examines whether both vectors are aligned in the same direction. This characteristic makes cosine similarity particularly suitable for gauging significant similarities between images. The formula for cosine similarity is provided below. In the given equations, T denotes the latent vector of the true image, whereas G represents the latent vector of the generated image.

$$\text{Cosine similarity} = 1 - \frac{T \times G}{\|T\| \|G\|} \quad (2)$$

The original images and the images generated through the data generating phase were passed to the autoencoder's encoder to obtain embeddings. The cosine similarity between the

100 original images and the generated image was calculated. The generated images with the top 10% of cosine similarity were selected and added to the training set.

C. Classification phase

To demonstrate the validity of our framework, we applied a classification model to images generated using our framework. The model we used for this purpose is YOLOv8, which is an enhancement of YOLOv5 based on additional layer modifications to improve the model's performance, achieving state-of-the-art results. The YOLO series is a well-known model extensively utilized in various CAD applications [12], [13].

The training dataset comprises a total of 10 sessions, divided into 5 sessions for the heart and 5 sessions for the kidney. We initiated the training with 100 samples and gradually increased the training dataset size to 500 samples. Since all data used to train the generation model are disease-related, we also labeled normal data to assess classification accuracy. Each class, including normal, was labeled with 500 samples per class. For testing, we extracted 50 data samples per class from the actual data that were not duplicated with the training data.

For evaluation metrics, we utilized the confusion matrix, along with accuracy, precision, recall, and F1-score [14], [15]. The confusion matrix is a table used in machine learning to assess the performance of a classification model, summarizing the relationship between the model's predictions and actual

TABLE I
THE RESULTS OF FID VALUES FOR GENERATING RADIOGRAPHIC IMAGE

Original		Cycle-1		Cycle-2		Cycle-3		Cycle-4	
Optimal	Worst	Optimal	Worst	Optimal	Worst	Optimal	Worst	Optimal	Worst
84.14	100.16	64.31	97.21	58.39	82.23	50.75	74.57	53.87	81.82

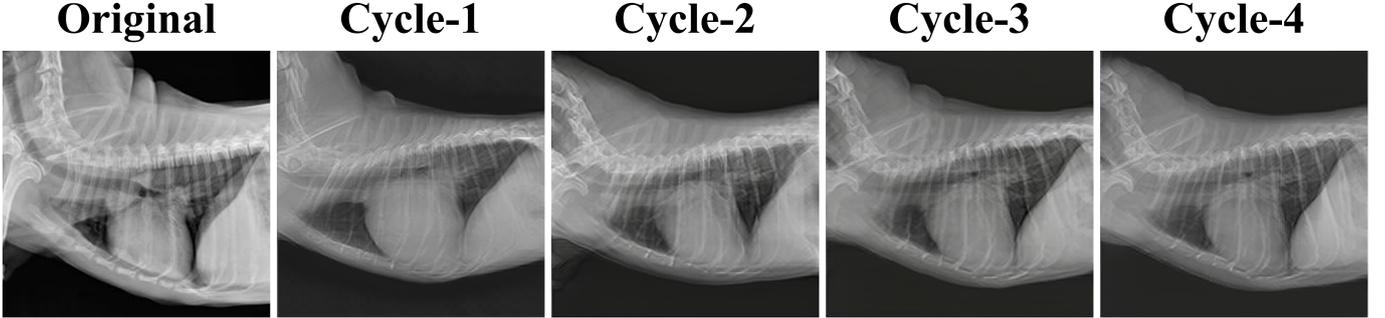


Fig. 3. Generation results of each cycle. The following is data generated by learning cardiomegaly data.

values. From this matrix, accuracy, precision, recall, and F1-score can be calculated using the following formulas. In these formulas, true positive (TP) represents the number of correctly predicted positive observations, while true negative (TN) denotes the number of correctly predicted negative observations. False positive (FP) indicates instances predicted as positive, but actually negative. False negative (FN) indicates the number of instances that are actually positive but are incorrect.

$$Accuracy : \frac{TP + TN}{TP + TN + FP + FN} \quad (3)$$

$$Precision : \frac{TP}{TP + FP} \quad (4)$$

$$Recall : \frac{TP}{TP + FN} \quad (5)$$

$$F1 \text{ score} : \frac{2 \cdot (Precision \cdot Recall)}{Precision + Recall} \quad (6)$$

III. EXPERIMENTS

A. Data generating phase

The results of the FID experiment to generate radiograph data are presented in Table I. In total, there are five sessions, ranging from original to cycle-4, each divided into ‘Optimal’ and ‘Worst’ cases. The original session involves training the model exclusively on selected data from the original dataset. ‘Optimal’ represents the lowest FID value among the 10 values, while ‘Worst’ indicates the highest FID value out of the 10. A lower FID value implies better performance of the generative model.

The term cycle refers to the process of generating data in the data generating phase, filtering through the query phase, and adding the filtered data to the training dataset of the generative model. The number following the cycle represents the iteration count. For example, cycle-1 involves training the generative model with 100 new data added to the training dataset, generated using the model trained on the data from the original session. Consequently, in cycle-1, the size of the training dataset is 200. Subsequently, cycle-2 includes the 200 data previously used and an additional 100 generated data.

In cycle-3, the generated radiograph data showed the best performance with the ‘Optimal’ score of 50.75 and the ‘Worst’ performance at 74.57. Additionally, a trend was observed in which performance tended to be less favorable with a smaller amount of data. On the contrary, as the amount of data increased, there was a performance improvement, although the final cycle did not consistently yield the best results for radiographic images. The generated data examples are shown in Fig. 3, where the left side displays the original data and the results obtained through cycle-4. In Fig. 3, the displayed results represent the filtered data for each cycle with the highest cosine similarity.

B. Classification phase

The results of the classification using YOLOv8 are presented in Fig. 4. Fig. 4 includes the confusion matrix, where the upper part represents the results tested on radiographs of dogs with cardiomegaly. The confusion matrix is commonly employed as an evaluation metric in various research studies involving classification tasks [16]. Each classification task was conducted to demonstrate the validity of the data and the classification model training utilized the dataset used for the generative model.

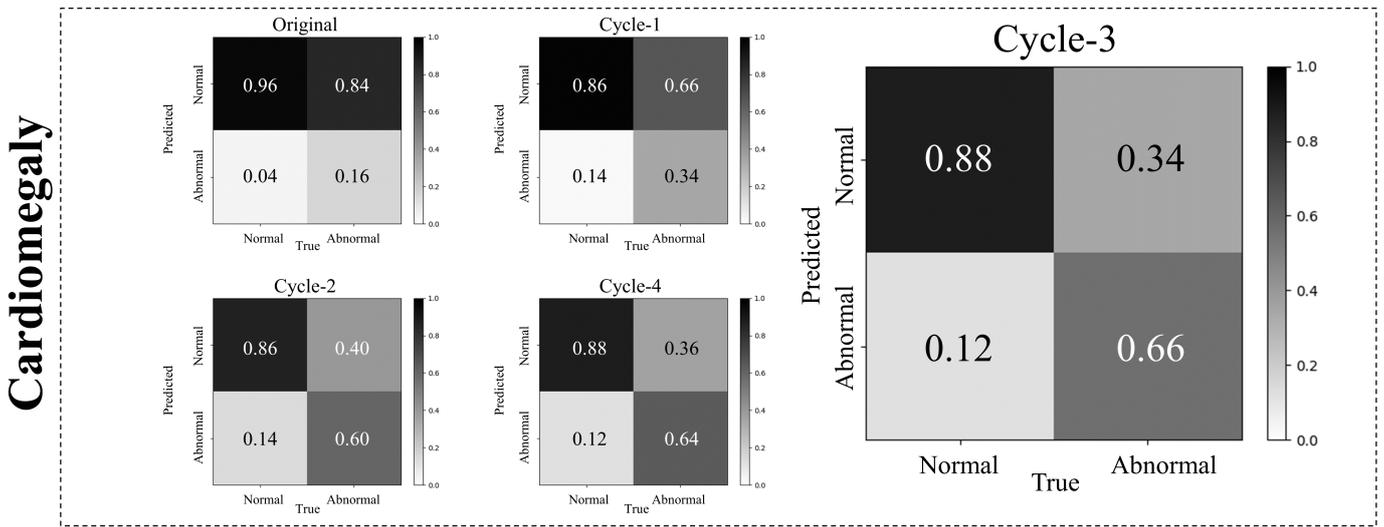


Fig. 4. Confusion matrix results of classification phase. The above results represent training and testing results using cardiomegaly data. Among the five sessions, the confusion matrix of the session with the highest accuracy is presented more prominently.

In cardiomegaly data, cycle-3 showed the highest accuracy and precision, with values of 0.73 and 0.72, respectively. In addition, cycle-3 and cycle-4 showed the highest F1-score of 0.79. The recall value was highest when trained with the original data.

On the other hand, in the case of cardiomegaly, the session with the lowest accuracy was not the original, but cycle-1. In cycle-1, the accuracy was 0.66, which was 0.01 lower than the accuracy of 0.67 in the original session.

IV. DISCUSSION

In this paper, we propose a generative active learning framework that automates the query process using the VAE. This framework generates data during the data generating phase and incrementally augments the training dataset of the generative model by filtering data through the query phase. Unlike previous research, we adopt the VAE to enhance the robustness of the query process. The query process has the filtering step by calculating the cosine similarity between the generated images and real images using 10% of the generated data. Experimental results demonstrate that iterative repetition of this process leads to improved performance of the generative model.

Observing the change in FID during the data generating phase, there was a consistent trend of FID reduction as the cycles progressed. For each session, 10 FID scores were obtained. In the case of radiographic data generation, the ‘Optimal’ FID score decreased from 84.14 to 50.75, and the ‘Worst’ FID score decreased from 100.16 to 74.57 throughout the cycles. These results demonstrate that our proposed framework effectively enhances the robustness of the generative model’s performance.

To validate the validity of our data, we conducted a classification phase. Data used for generation belonged to all cate-

gories of disease, including cardiomegaly. The experimental results, as observed from the confusion matrix, reveal that as cycles increase, the increase in disease data leads to an increase in FP. In the case of cardiomegaly, FP increased from 0.16 to 0.66. However, the overall accuracy for both diseases reached its highest value at cycle-3. It is important to note that when evaluating the model’s performance numerically, the classification performance of normal data also influences the results. Therefore, while accuracy has increased slightly, the significant increase in FP suggests that the generated data using our framework has demonstrated its utility in improving the performance of the classification model.

Furthermore, this study has some limitations. First, our methodology involved the utilization of an existing GAN variant instead of the proposed model. In particular, contemporary image generation models lean toward diffusion models [17], [18] rather than GAN. Second, while there is a plethora of medical image data available, we restricted our application of the framework to radiographic image data. Given the limited dataset that encompasses only these two modalities, further exploration is essential across diverse datasets to establish the generalizability of our findings. These limitations highlight avenues for future research and improvements in our approach.

V. CONCLUSION AND FUTURE WORKS

This study proposed the VAE-based generative active learning framework. The potential of this framework to address the issue of medical data scarcity in CAD was demonstrated through experimental results, including the FID of the generative model and the confusion matrix, accuracy, F1-score, precision, and recall of the classification model. Future research will extend to proposing generative model such as diffusion models and using various types of data, such as computerized tomography (CT), magnetic resonance imaging (MRI), etc. It

will have a positive impact on the performance improvement of the CAD system in the future and provide an opportunity to promote the development of the medical AI field.

REFERENCES

- [1] X. Zhu, C. Vondrick, C. C. Fowlkes, and D. Ramanan, "Do we need more training data?," *International Journal of Computer Vision (IJCV)*, vol. 119, no. 1, pp. 76–92, 2016.
- [2] V. Gudivada, A. Apon, and J. Ding, "Data quality considerations for big data and machine learning: going beyond data cleaning and transformations," *International Journal on Advances in Software*, vol. 10, no. 1, pp. 1–20, 2017.
- [3] J.-Y. Oh, I.-G. Lee, H.-H. Chang, E. Lee, and J.-H. Jeong, "Application of a dual-stage deep learning framework to detect left atrial enlargement for pet heart failure," *IEEE International Conference on Systems, Man, and Cybernetics (SMC), Hawaii, USA, Oct. 1-4, 2023*.
- [4] D. Yoon, H.-J. Kong, B. S. Kim, W. S. Cho, J. C. Lee, M. Cho, M. H. Lim, S. Y. Yang, S. H. Lim, and J. Lee, "Colonoscopic image synthesis with generative adversarial network for enhanced detection of sessile serrated lesions using convolutional neural network," *Scientific Reports*, vol. 12, no. 1, p. 261, 2022.
- [5] M. La Salvia, E. Torti, R. Leon, H. Fabelo, S. Ortega, B. Martinez-Vega, G. M. Callico, and F. Leporati, "Deep convolutional generative adversarial networks to enhance artificial intelligence in healthcare: a skin cancer application," *Sensors*, vol. 22, no. 16, p. 6145, 2022.
- [6] J.-J. Zhu and J. Bento, "Generative adversarial active learning," *arXiv preprint arXiv:1702.07956*, 2017.
- [7] D. P. Kingma and M. Welling, "Auto-encoding variational bayes," *arXiv preprint arXiv:1312.6114*, 2013.
- [8] C. Lam, B. J. Gavaghan, and F. E. Meyers, "Radiographic quantification of left atrial size in dogs with myxomatous mitral valve disease," *Journal of Veterinary Internal Medicine (JVIM)*, vol. 35, no. 2, pp. 747–754, 2021.
- [9] A. Sauer, K. Chitta, J. Müller, and A. Geiger, "Projected GANs converge faster," *Advances in Neural Information Processing Systems (NIPS)*, vol. 34, pp. 17 480–17 492, 2021.
- [10] M. Heusel, H. Ramsauer, T. Unterthiner, B. Nessler, and S. Hochreiter, "Gans trained by a two time-scale update rule converge to a local nash equilibrium," *Advances in Neural Information Processing Systems (NIPS)*, vol. 30, 2017.
- [11] F. Rahutomo, T. Kitasuka, and M. Aritsugi, "Semantic cosine similarity," in *International Student Conference on Advanced Science and Technology (ICAST), Seoul, South Korea, Oct. 29-30*, vol. 4, no. 1, 2012, p. 1.
- [12] A. K. Chaudhary, S. Roy, R. Rizk, and K. Santosh, "Automated fracture detection from CT scans," in *IEEE International Conference on Artificial Intelligence (CAI)*, 2023, pp. 161–162.
- [13] J. Estrada, Y. Zhigang, S. Datta, N. Duraisamy, J. De Guia, O. Cheng Hun, G. Opina, and A. Tripathi, "AV in action: a development of robust and efficient planning and perception system for autonomous food delivery vehicle," in *IEEE International Conference on Artificial Intelligence (CAI)*, 2023, pp. 19–20.
- [14] D. M. Powers, "Evaluation: from precision, recall and F-measure to ROC, informedness, markedness and correlation," *arXiv preprint arXiv:2010.16061*, 2020.
- [15] A. Tharwat, "Classification assessment methods," *Applied Computing and Informatics*, vol. 17, no. 1, pp. 168–192, 2020.
- [16] J.-H. Jeong, J.-H. Cho, B.-H. Lee, and S.-W. Lee, "Real-time deep neuro-linguistic learning enhances noninvasive neural language decoding for brain-machine interaction," *IEEE Transactions on Cybernetics*, vol. 53, no. 12, pp. 7469–7482, 2023.
- [17] J. Ho, A. Jain, and P. Abbeel, "Denoising diffusion probabilistic models," *Advances in Neural Information Processing Systems (NIPS)*, vol. 33, pp. 6840–6851, 2020.
- [18] Song, Jiaming and Meng, Chenlin and Ermon, Stefano, "Denoising diffusion implicit models," *arXiv preprint arXiv:2010.02502*, 2020.