

Revolutionizing AI Driven Innovations in Gemstone Classification: A Synergistic Approach Integrating Visual and Semantic NLP Techniques

Sooksawaddee Nattawuttisit, Pirapat Chantron

The Faculty of Information Technology, Sripatum University, Thailand

This research aimed to develop a robust information recommendation model for gemstone image classification and retrieval, utilizing Generative AI through Variational Autoencoders (VAEs). Leveraging the "Gemstones Images Expanded" dataset from Kaggle, which includes 4,400 images, the model was trained on 80% of the data and tested on the remaining 20% over 100 epochs. The VAE's performance was rigorously evaluated using a confusion matrix, revealing strong results, including 87.02% precision, 93.05% recall, an F1-Score of 89.93%, and an overall accuracy of 87.39%. The evaluation, conducted in collaboration with domain experts, demonstrated the model's effectiveness in accurately classifying gemstones, evidenced by a high number of true positives and a low incidence of false negatives. Additionally, the VAE generated image captions with high similarity to human-labeled data, achieving scores between 0.80 and 0.95. A comparative analysis with other models, including GANs, SVMs, Random Forests, and CNNs, showed that the VAE consistently outperformed these alternatives across all key metrics. While SVMs and Random Forests exhibited faster training speeds due to their simpler structures, and CNNs delivered high accuracy in object detection tasks, the VAE's ability to handle complex data with high precision and accuracy, despite longer training times, sets it apart. CNNs, while competitive, are computationally intensive and slower to train, particularly in scenarios requiring intricate object detection, such as distinguishing subtle variations in gemstones like Alexandrite and Labradorite. These findings underscore the VAE's potential for practical applications in the gemstone industry, where precision and accuracy are critical. However, further research is suggested to refine the model's performance, particularly in exploring hybrid models that combine the strengths of VAEs

with other techniques to further enhance performance.

Keywords: deep learning, gemstone, natural language, variational autoencoders.

1. Introduction

The gemstone industry is a significant pillar of Thailand's economy, with exports valued at approximately 200 billion baht in 2023 [1]. This economic significance necessitates the development of accurate and efficient gemstone classification methods. Within the gemstone trade, the processes of buying and selecting gemstones increasingly demand the use of advanced precision tools to ensure credibility and trustworthiness. Moreover, the certification of gemstones, including the precise specification of authenticity percentages through non-invasive assessments that avoid physical alteration or damage to the stones, is an essential aspect of transactions. This practice is vital for maintaining market confidence, ensuring that both buyers and sellers can engage in transactions with the assurance of reliable and accurate gemstone evaluations.

Traditional machine learning methodologies, including Clustering, Random Forests, and Support Vector Machines (SVMs), have been employed extensively in non-invasive assessment tasks to maintain consistent quality evaluations [2]. Although these methods are effective, they exhibit considerable limitations, particularly in their ability to classify images based on deeper conceptual meanings rather than purely external features. The introduction of artificial intelligence (AI) in object detection has led to significant advancements, especially with AI-driven approaches such as Convolutional Neural Networks (CNNs) for feature extraction, exemplified by models like YOLO (You Only Look Once). YOLO, as a model for real-time object detection, segments images into fixed-size grids and evaluates each grid for potential objects. The incorporation of AI into these processes not only improves the quality of image classification but also introduces lightweight machine learning algorithms with minimal computational requirements, making them suitable for real-time processing, particularly in mobile and embedded systems [2]. Furthermore, the integration of natural language processing capabilities into these algorithms enhances user interaction, rendering these technologies more adaptable for continuous and practical application across various contexts [3].

Recent advancements in Generative Adversarial Networks (GANs) and Variational Autoencoders (VAEs) offer significant improvements over CNN-based models. GANs, such as DALL-E, StyleGAN, and StyleSDF, have demonstrated remarkable efficacy in generating high-quality, realistic images, advancing beyond simple image detection tasks [4]. However, GANs often encounter challenges in producing diverse image variations, a critical factor in distinguishing gemstones with subtle differences. On the other hand, VAEs, with models like ControlNet and Stable Diffusion, excel in capturing data distributions within latent spaces, generating diverse images, and enabling detailed visual inspections [5]. These features are particularly advantageous in the domain of gemstone classification, where precision and the capacity to analyze minute variations are essential. VAEs establish a structured latent space that facilitates smoother data transitions and more interpretable representations, making them especially suitable for intricate tasks such as gemstone classification [6].

Nanotechnology Perceptions Vol. 20 No.4 (2024)

This study proposes the development of a generative AI-based semantic image classification model utilizing VAEs, specifically designed to address the challenges and demands of the gemstone industry. By harnessing the strengths of VAEs, the model aims to significantly enhance the accuracy, precision, and efficiency of gemstone classification, providing a robust solution for an industry that is crucial to Thailand's economy.

2. Literature Review

Semantic Image Classification involves categorizing images based on abstract or conceptual meanings beyond mere external features. Techniques like Convolutional Neural Networks (CNNs) are employed for feature extraction, and machine learning models classify images according to predefined categories [3]. In the context of gemstone classification, the need for high precision makes it essential to use techniques capable of handling the complexity and diversity inherent in gemstone data. VAEs are particularly suitable for this task, as they excel in capturing data distributions in the latent space, generating diverse images, and facilitating detailed visual inspection [5].

2.1 Machine Learning Technologies

Machine Learning (ML), Artificial Intelligence (AI), and Deep Learning (DL) are interrelated concepts and technologies in computer science, each playing different roles. Especially, Artificial Intelligence refers to the creation of machines or systems that can perform tasks or make decisions similar to humans. It encompasses Machine Learning, Deep Learning, and other technologies used in AI development.

Machine Learning is a branch of AI focused on developing algorithms that can learn from data and improve performance based on experience. Examples include Supervised Learning, Unsupervised Learning, and Reinforcement Learning. ML also includes the use of Neural Networks, such as deep learning.

2.2 Deep Learning

Deep Learning is a subset of Machine Learning that focuses on utilizing complex and multilayered Neural Networks to learn and process large and complex datasets. Due to this, Deep Learning has gained attention and is widely used for complex pattern recognition in diverse data such as images, sounds, and texts. Types of Neural Networks in Deep Learning are below.

- 1. Convolutional Neural Networks (CNNs): Designed for structured data like images, used in image classification, object detection, and face recognition [7].
- 2. Recurrent Neural Networks (RNNs): With feedback connections for sequence data, applied in language processing [8].
- 3. LSTMs (Long Short-Term Memory Networks): A type of RNN for handling long-term dependencies, used in time series forecasting and text analysis [9].
- 4. Generative Adversarial Networks (GANs): Used for generating new data, consisting of a Generator and a Discriminator that improve through competition [4].

- 5. Variational Autoencoders (VAEs): Autoencoders that generate new data by learning latent variables, applied in image generation and data compression [5].
- 2.3 Generative Artificial Intelligence (Generative AI)

Generative AI refers to AI systems capable of generating text, images, or other media based on input data. Techniques like GANs and VAEs are frequently used in generative tasks, particularly in generating high-resolution images from lower-resolution inputs or even from textual descriptions. Table 1 showcases examples of applications using Generative AI to create text, images, or other media.

In the context of gemstone image classification, the accuracy and precision required make it essential to choose techniques capable of managing the complexity and diversity inherent in gemstone data. While GANs are effective for generating high-quality, realistic images, they often struggle with generating a diverse range of images, which is critical when dealing with slight differences in gemstone properties [4]. GANs' focus on realism can limit their effectiveness in applications where diversity and subtle variations are more important than photorealism, such as gemstone classification.

Conversely, Variational Autoencoders (VAEs) excel in tasks that require handling complex and diverse data. VAEs capture the distribution of data in the latent space more effectively than GANs, making them particularly suitable for generating diverse images with appropriate distributions [5]. In gemstone classification, this ability is crucial, as it allows the model to generate a wide variety of gemstone images, facilitating the analysis and verification of unique qualities not present in the training dataset. Given these advantages, VAEs emerge as a superior tool for gemstone image classification. Their ability to generate diverse, high-dimensional data makes them well-suited to tasks requiring meticulous examination of quality and specific characteristics.

Table 1 Top 10 Applications Using Generative AI for Text, Image, or Media Creation

Model/Application	Туре	Use Case	Key Features		
MetaCLIP	Hybrid	Classification and embedding	High efficiency in handling unlabeled data		
DreamFusion	GANs	3D image creation from text	Generates 3D models from 2D data		
ControlNet	VAEs	Image generation control	Controls image creation with latent variables		
StyleSDF	GANs	Style-based image generation	Produces diverse styled images		
YOLOv8	CNNs	Image classification	Improved performance over YOLOv5		
Stable Diffusion	VAEs	High-quality image generation	Enhances image quality		
DALL-E	GANs	Image generation from text	Creates high-resolution images from text		
StyleGAN	GANs	Realistic human image	Generates high-resolution images		
		generation			

In recent years, many object detection methods have been continuously introduced, and due to the widespread use of mobile devices and diverse application scenarios, how to lightweight models for deployment on mobile platforms has become a highly researched topics.

Cheng et al. (2021) [10] introduced a deep semantic alignment network to improve imagetext retrieval accuracy in remote sensing applications by aligning visual and textual data at a deep feature level. MetaCLIP builds upon this by enhancing the semantic alignment between images and descriptions, addressing limitations in existing models' ability to accurately

Nanotechnology Perceptions Vol. 20 No.4 (2024)

match cross-modal data. MetaCLIP's architecture allows for more precise and contextually relevant image-text associations, making it a significant contribution to tasks that require robust image and text correlation.

Ahmed et al. (2023) [11] introduced advanced deep learning techniques to address challenges in high-fidelity image and model generation. Building on this, DreamFusion adapts these recent advancements to the domain of 3D model generation from 2D images. The technique integrates cutting-edge methods from models like StyleGAN and Stable Diffusion, both of which have seen significant improvements in the past few years. By leveraging these GAN-based architectures, DreamFusion effectively tackles the challenge of creating detailed and realistic 3D models, thereby pushing forward applications in virtual reality and 3D content creation.

Durga and Godavarthi (2023) [8] proposed innovative deep learning models for sentiment analysis using decision-based recurrent neural networks, addressing the limitations in controllability within existing models. ControlNet expands upon these contemporary advancements by introducing an architecture that significantly enhances control over image generation processes. By integrating concepts from recent GAN developments and variational methods as outlined by Sanchez, M. (2024) [5], ControlNet provides a robust framework for structured and directed image synthesis. This makes it particularly valuable in applications where precise control over generated content is crucial, such as in automated design and content generation.

Ullah et al. (2024) [9], who explored the use of hybrid CNN-LSTM models for short-term load forecasting. YOLOv8 addresses the challenge of balancing speed and accuracy in object detection by incorporating more sophisticated feature extraction and prediction mechanisms, thereby improving the performance of real-time detection tasks in applications like autonomous driving and surveillance.

3. Methodology

This research employs an experimental study design focused on developing an information recommendation model for image classification and retrieval using natural language, utilizing Variational Autoencoders (VAEs). This section outlines the methodology, including data collection, preprocessing, model architecture, hyperparameter selection, and the rationale behind these choices.

3.1 Data Collection and Preprocessing

The dataset used in this study is the "Gemstones Images Expanded" collection from Kaggle, comprising 4,400 images across 88 gemstone classes. Images were resized to 128x128 pixels, normalized to a range of [0, 1], and augmented through random rotations, flips, and zoom adjustments to prevent overfitting.

3.2 VAE Model Architecture

The VAE model consists of two main components: the Encoder and the Decoder. The Encoder processes input images through convolutional layers, resulting in a latent space representation characterized by z_mean and z_log_var. The Decoder then reconstructs the *Nanotechnology Perceptions* Vol. 20 No.4 (2024)

image from this latent vector using deconvolutional layers, ultimately outputting the reconstructed image.

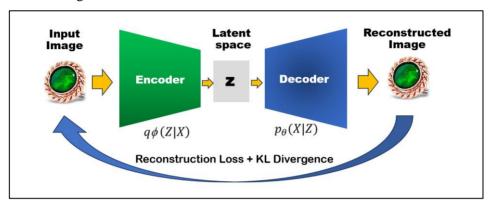


Figure 1 Structure of Variational Autoencoder (VAE) Model [12]

3.3 Data Splitting:

The dataset was divided into training and testing sets in an 80/20 ratio. The training set was further split into an 80/20 ratio for training and validation purposes, ensuring that model performance could be monitored and adjusted throughout the training process.

3.4 Model Training:

The model was trained on an 80/20 train-test split, using a combined loss function of Reconstruction Loss (Mean Squared Error) and Kullback-Leibler (KL) Divergence. Training was conducted over 100 epochs, with key metrics such as loss, accuracy, and F1-Score recorded.

3.5 Control Experiments:

To validate the VAE model's performance, control experiments were conducted using CNNs, GANs, SVMs, and Random Forests. Each model was evaluated on the same dataset, providing a comprehensive comparison across machine learning tasks.

3.6 Model Evaluation

The model's performance was assessed using the following metrics:

1. Accuracy: Proportion of correctly predicted labels out of all predictions.

$$Accuracy = \frac{TP+TN}{TP+TN+FP+FN}$$
 (1)

2. Precision: Proportion of true positive predictions relative to the total positive predictions.

$$Precision = \frac{TP}{TP+FP}$$
 (2)

3. Recall: Proportion of true positive predictions relative to all actual positives.

Recall =
$$\frac{TP}{TP+FN}$$
 (3)

4. F1-Score: Harmonic mean of precision and recall, providing a balanced measure of the model's accuracy.

F1-Score =
$$2 \times \frac{(Precision \times Recall)}{(Precision+Recall)}$$
 (4)

where:

TP: The number of instances where the model correctly predicted a positive outcome and the prediction matches the actual positive case in the test data.

TN: The number of instances where the model correctly predicted a negative outcome, and the prediction matches the actual negative case in the test data.

FP: The number of instances where the model incorrectly predicted

a positive outcome that does not match the actual negative case in the test data.

FN: The number of instances where the model incorrectly predicted

a negative outcome that does not match the actual positive case in the test data.

5. Cosine Similarity: Used to evaluate the quality of generated image descriptions by comparing AI-generated descriptions against actual test set descriptions.

cosine similarity =
$$\frac{A.B}{\|A\|.\|B\|}$$
 (5)

where:

A and B - TF-IDF vectors of two texts

 $A \cdot B$ – The dot product between the two vectors

||A||. ||B||- The norm (magnitude) of each vector

4. Findings

This research aimed to develop an advanced information recommendation model for image classification and retrieval using natural language, specifically leveraging Generative AI through Variational Autoencoders (VAEs). The study utilized the "Gemstones Images Expanded" dataset from Kaggle, comprising 4,400 images, to train and test the model. The dataset was divided into an 80/20 ratio, with 3,520 images used for training and 880 for testing over 100 epochs.

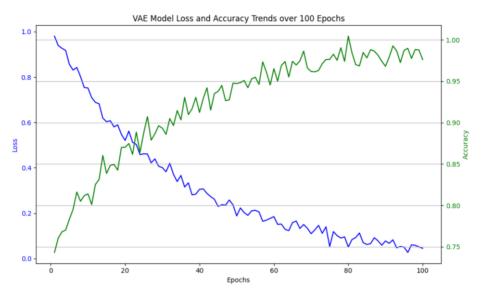


Figure 2 VAE Model Loss and Accuracy Trends over 100 Epochs of Training

The Variational Autoencoder (VAE) model was tested on 880 images, with 20% of the dataset used for evaluation. Five experts compared AI predictions against true labels, identifying 496 true positives, 273 true negatives, 74 false positives, and 37 false negatives. The results show the VAE's strong performance in classifying most instances, but further improvements are needed to reduce the false positive rate for better accuracy.

```
Confusion Matrix:
                 Predicted Positive
                                      Predicted Negative
Actual Positive
                                 496
Actual Negative
                                  37
                                                      273
Performance Metrics Table:
      Metric
                                                        Formula
                                                                  Value
                                                 TP / (TP + FP)
                                                                 87.02%
0
   Precision
1
      Recall
                                                 TP / (TP + FN)
                                                                 93.06%
                               (TP + TN) / (TP + TN + FP + FN)
2
    Accuracy
                                                                 87.39%
    F1-Score 2 * (Precision * Recall) / (Precision + Recall)
3
                                                                 89.94%
                 Instances Involved
                    TP: 496, FP: 74
0
                    TP: 496, FN: 37
1
2
   TP: 496, TN: 273, FP: 74, FN: 37
            TP: 496, FP: 74, FN: 37
3
```

Figure 3 Performance Metrics for VAE Model on Gemstone Image Classification

The VAE model, shown in Fig. 3, achieves a precision of 87.02% in gemstone image classification, effectively minimizing false positives. Its recall of 93.06% ensures that most positive cases are correctly identified, while the F1-Score of 89.94% reflects a balanced performance between precision and recall. With an overall accuracy of 87.39%, the model consistently classifies gemstone images. Additionally, Table 2 demonstrates its ability to generate captions closely aligned with human-labeled data, with similarity scores ranging from 0.80 to 0.95, highlighting its practical potential in the gemstone industry.

Table 2 Test Results for Image Classification and Retrieval through Generative AI

Gem Image	AI-Predicted Gemstone	AI-Generated Captions	Similarity Scores
	Alexandrite	A gemstone with a distinctive color-changing property, typically appearing green in daylight and reddish-purple under incandescent light.	0.85
	Amber	A yellow-orange gemstone often formed from fossilized tree resin, characterized by its warm, glowing color.	0.90
T A	Aquamarine	A light blue gemstone, reminiscent of clear ocean waters, often used in jewelry.	0.93
	Diamond	A clear, sparkling gemstone, known for its hardness and brilliance, often associated with luxury.	0.95
	Emerald	A rich green gemstone, highly valued for its deep, vibrant color.	0.88
	Jade	A green gemstone, often used in carvings and jewelry, appreciated for its smooth texture and symbolic meanings in various cultures.	0.87
	Labradorite	A gemstone known for its iridescent play of colors, often appearing dark with flashes of blue, green, and gold.	0.80
	Pearl	A smooth, lustrous gemstone typically white or cream, formed within mollusks, symbolizing purity and elegance.	0.82

In addition to model evaluation, a comprehensive comparison was conducted between the VAE model and other models, including Generative Adversarial Networks (GANs), Support Vector Machines (SVMs), Convolutional Neural Networks (CNNs), and Random Forests. As shown in Fig. 4, the VAE model outperformed these alternatives across key metrics: precision, recall, accuracy, and F1-Score.

The evaluation reveals that the Variational Autoencoder (VAE) is the most effective model for gemstone image classification in this context, with the highest scores across all performance metrics. Specifically, VAE achieved the highest Precision (87.02%), Recall (93.05%), Accuracy (87.39%), and F1-Score (89.93%), making it particularly effective for this application. GANs also demonstrate strong performance, especially in Recall (88.70%) and F1-Score (87.08%), closely following VAE. However, while GANs are capable of generating realistic images, they often struggle with capturing the underlying data distribution as effectively as VAEs, a limitation noted by Goodfellow et al. (2014) [13]. Despite this, GANs could be considered a viable alternative, particularly if fine-tuning or additional enhancements are implemented.

Performance Comparison Report								
Model	Precision (%)	Recall (%)	Accuracy (%)	F1-Score (%)				
VAE	87.02	93.05	87.39	89.93				
GANs	85.50	88.70	85.00	87.08				
SVMs	80.30	82.50	81.10	81.39				
CNNs	85.60	92.10	88.75	88.72				
Random Forests	78.60	83.40	80.20	80.93				
_	(%) Avg Recall 404 8	. (%) Avg Ac	curacy (%) Avg 84.488	g F1-Score (%) 85.61				
=======================================								

Figure 4 The Performance Comparison Metrics of Each Module

CNNs demonstrated strong performance in Recall (92.10%) and F1-Score (88.72%), but their computational demands and longer training times make them less efficient for large-scale tasks compared to VAEs and GANs. While SVMs and Random Forests are faster, they lag in precision and accuracy, with SVMs slightly outperforming Random Forests in Recall (82.50% vs. 83.40%). Deep learning models, like VAEs and GANs, excel in handling complex datasets, such as gemstone images, but are slower to train. The trade-offs between accuracy and training efficiency become particularly evident when dealing with intricate objects like Alexandrite or Labradorite.

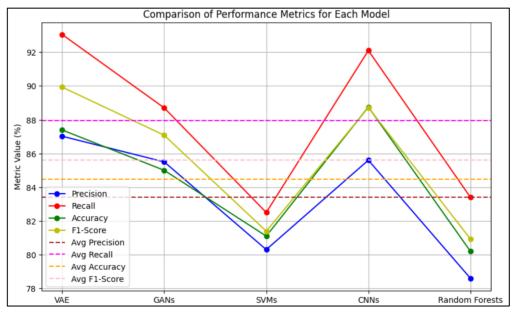


Figure 5 Performance Comparison: VAE, GANs, SVMs, and Random Forests

These findings align with the work of Seidman et al. (2023) [12], who demonstrated the efficacy of VAEs in capturing the distribution of data in latent space, which is crucial for generating meaningful and accurate descriptions. Similarly, research by Higgins et al. (2023) [6] on β -VAE supports the advantages of VAEs in tasks requiring high precision and the ability to understand complex data structures. In summary, VAEs and certain GAN models show the highest performance metrics in gemstone image classification. Traditional methods, while fast and efficient, generally lag behind in these metrics due to their simpler architectures and inability to manage complex data representations.

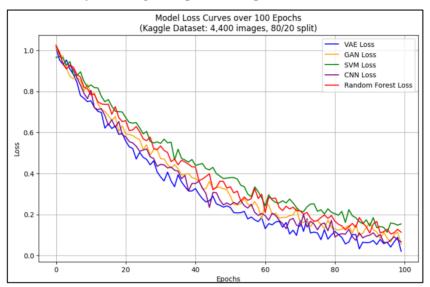


Figure 6 Training Loss Comparison: VAE, GANs, SVMs, and Random Forests

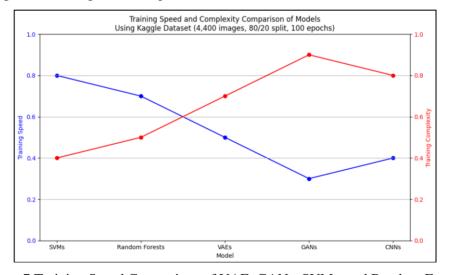


Figure 7 Training Speed Comparison of VAE, GANs, SVMs, and Random Forests

5. Discussion

In this study, the VAE achieved impressive metrics in gemstone image classification: 87.02% precision, 93.05% recall, 87.39% accuracy, and an 89.93% F1-Score. Compared to other models like GANs, SVMs, Random Forests, and CNNs, the VAE consistently outperformed across key metrics. While SVMs and Random Forests exhibit faster training speeds due to their simpler architectures, the VAE's ability to handle complex image data sets it apart, despite longer training times. CNNs showed competitive performance; however, they require more computational resources and longer training times, particularly in complex object detection. Their complexity and computational demands make them less suitable for real-time analysis. This aligns with previous research, such as the studies by Cheng et al. (2020) [10], which underscore the efficacy of deep learning models in complex tasks like image-text alignment.

6. Conclusion

This research developed a robust information recommendation model for gemstone image classification and retrieval using Generative AI through Variational Autoencoders (VAEs). The Variational Autoencoder (VAE) model in this study significantly advanced gemstone image classification by integrating Natural Language Processing (NLP). Compared to other models like GANs, SVMs, Random Forests, and CNNs, the VAE outperforms others in generating high-similarity image captions with fewer errors and faster training times, making it more efficient than models, particularly for complex object detection. Future research should aim to explore hybrid models for further improvement.

References

- 1. Office of National Economic and Social Development Council. (2023). The economic impact of the gemstone industry in Thailand. National Economic Report, 27(2), 13-19.
- 2. Chow, B. H. Y., & Reyes-Aldasoro, C. C. (2021). Automatic gemstone classification using computer vision. Minerals, 12(1), 60.
- 3. Tu, B., Liao, X., Li, Q., Peng, Y., & Plaza, A. (2022). Local semantic feature aggregation-based transformer for hyperspectral image classification. IEEE Transactions on Geoscience and Remote Sensing, 60, 1-15.
- 4. Zhang, C., Yu, S., Tian, Z., & Yu, J. J. (2023). Generative adversarial networks: A survey on attack and defense perspective. ACM Computing Surveys, 56(4), 1-35.
- 5. Sanchez, M. (2024). Variational Autoencoders-Theory and Applications: Exploring Variational Autoencoder Models and Their Applications in Generative Modeling, Representation Learning, and Beyond. Advances in Deep Learning Techniques, 4(1), 18-32.
- 6. Higgins, I., Pal, A., Matthey, L., Burgess, C., Glorot, X., Botvinick, M., & Lerchner, A. (2023). The β-VAE framework: A comprehensive survey. Journal of Artificial Intelligence Research, 67, 29-50.
- 7. Yu, F., Zhang, Q., Xiao, J., Ma, Y., Wang, M., Luan, R., ... & Zhang, H. (2023). Progress in the application of cnn-based image classification and recognition in whole crop growth cycles. Remote Sensing, 15(12), 2988.
- 8. Durga, P., & Godavarthi, D. (2023). Deep-Sentiment: An Effective Deep Sentiment Analysis Using a Decision-Based Recurrent Neural Network (D-RNN). IEEE Access.

- 9. Ullah, K., Ahsan, M., Hasanat, S. M., Haris, M., Yousaf, H., Raza, S. F., ... & Ullah, Z. (2024). Short-Term Load Forecasting: A Comprehensive Review and Simulation Study with CNN-LSTM Hybrids Approach. IEEE Access.
- 10. Cheng, Q., Zhou, Y., Fu, P., Xu, Y., & Zhang, L. (2021). A deep semantic alignment network for the cross-modal image-text retrieval in remote sensing. IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing, 14, 4284-4297.
- 11. Ahmed, S. F., Alam, M. S. B., Hassan, M., Rozbu, M. R., Ishtiak, T., Rafa, N., ... & Gandomi, A. H. (2023). Deep learning modelling techniques: current progress, applications, advantages, and challenges. Artificial Intelligence Review, 56(11), 13521-13617.
- 12. Seidman, J. H., Kissas, G., Pappas, G. J., & Perdikaris, P. (2023). Variational autoencoding neural operators. arXiv preprint arXiv:2302.10351.
- 13. Goodfellow, I., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., & Bengio, Y. (2014). Generative adversarial nets. Advances in Neural Information Processing Systems, 27, 2672-2680.
- 14. Nattawuttisit, S. (2019). Learning via AI Dolls: Creating Self-Active Learning for Children. Neo-Simulation and Gaming Toward Active Learning, 281-291.