

Enhanced image steganalysis through reinforcement learning, generative adversarial networks, and vision transformers (RGV-Stega)

Sumia Abdulhussien Razooqi Al-obaidi ¹, Mohammed Ahmed Talab ², Muhanaad Shakir ³,
Mustafa Talal Alnaseri ⁴, Suryanti Awang ^{5*}

¹Ministry of Higher Education and Scientific Research, Baghdad, Iraq

²Department of Medical Physics, College of Applied Science, Al-Fallujah University, Anbar, Iraq

³College of Business(COB), University of Buraimi, Oman

⁴Faculty of Electrical and Computer Engineering, University of Tabriz, Tabriz, Iran

⁵Faculty of Computing, Universiti Malaysia Pahang Al-Sultan Abdullah, Pekan, Pahang 26600, Malaysia

ARTICLE INFO

Received: 22/09/2025
Accepted: 18/12/2025
Available online: 24/03/2026
April Issue
[10.37652/juaps.2025.165053.1707](https://doi.org/10.37652/juaps.2025.165053.1707)

 CITE @ JUAPS

Corresponding author
Suryanti Awang
suryanti@umpsa.edu.my

ABSTRACT

As steganographic techniques advance and object data are increasingly generated in various domains, detecting concealed data embedded in digital media has become more challenging; thus, image steganalysis plays an important role. Conventional detection techniques are becoming less effective against these evolving strategies, which calls for better solutions. In this paper, we present RGV-Stega, a new high-level structure that combines Reinforcement Learning (RL), Generative Adversarial Networks (GANs), and Vision Transformers (ViTs) to further enhance the detection of steganographic contents in images. Artificial intelligence techniques such as RL, GANs, and ViTs play significant roles in current steganalysis work: RL helps an agent learn optimal feature extraction strategies, GANs yield rich steganographic samples for strong training, and ViTs help capture long-range dependencies in image data, thereby boosting top-1 accuracy. On benchmark datasets such as BOSSBase and BOWS, our method RGV-Stega achieves up to 93% accuracy, comparable to conventional CNN- and ViT-based models. Findings show that combining RL, GANs, and ViTs is an effective approach to addressing issues arising from advances in steganographic techniques. This improves steganalysis's ability to identify hidden content.

Keywords: *Generative adversarial network, Image steganalysis, Reinforcement learning, Transfer learning, Vision transformers.*

1 INTRODUCTION

As digital communication has evolved rapidly, it has become increasingly crucial to keep sensitive data confidential from unauthorized users. Encrypting messages in digital media is known as steganography. This approach allows people to communicate secretly. However, it can be harmful if used for malicious purposes. Steganalysis, the opposite of steganography, aims to uncover hidden messages. However, the increasing number of steganographic methods makes effective steganalysis

challenging [1].

Conventional steganalysis techniques often rely on manually constructed statistical features and traditional machine learning algorithms, which struggle to generalize across different embedding schemes and image domains [2]. Deep learning techniques have recently demonstrated promising progress across numerous computer vision tasks, including steganalysis. Recently, single-architecture models exhibit limitations in effectively capturing both spatial and frequency-domain

perturbations in stego images [3].

This paper proposes an integrated architecture that utilizes Reinforcement Learning (RL), Generative Adversarial Networks (GANs), and Vision Transformers (ViTs) to enhance the effectiveness of steganalysis, as shown in Figure 1. The paper presents RGV-Stega, a unified framework for steganographic analysis that incorporates Reinforcement Learning (RL), Generative Adversarial Networks (GANs), and Vision Transformers (ViTs) to address the increasing complexity of modern steganographic methods. The research constructs an adaptive, multi-view detection model capable of identifying both local and global visual perturbations. This indicates that conventional steganalysis methods, which have heavily relied on manually crafted features, struggle to distinguish high-fidelity embedding artifacts.

Our motivation stems from the complementary strengths of these paradigms:

- A GAN network, producing stego-cover pairs which are real and diverse, and which are used to reinforce training in diverse embedding conditions;
- An extractor based on ViT, which achieves long-distance global dependencies and finer-grained spatial interaction that CNNs generally do not have;
- A reward-based policy of adaptively choosing the most informative features and maximizing the accuracy of detection decisions can be found in an RL module.

Significant experimentation has been conducted on the BOSSBase and BOWS datasets, and our model achieves state-of-the-art accuracy (93.6%), surpassing benchmarks for CNN-, transformer-, and hybrid-driven steganalysis. Further, ablation research shows that each component provides incremental performance gains, and their combination yields the best results. Other analyses—such as reward shaping, choice of a loss function, and CNN depth—confirm the strengths of the model, and its sensitivity to the parameters of the model and generalization. We hypothesize that their synergistic combination will result in a robust steganalysis framework capable of detecting complex embedding strategies. The proposed method is effective under similar evaluations and similar datasets for different image augmentation attacks.

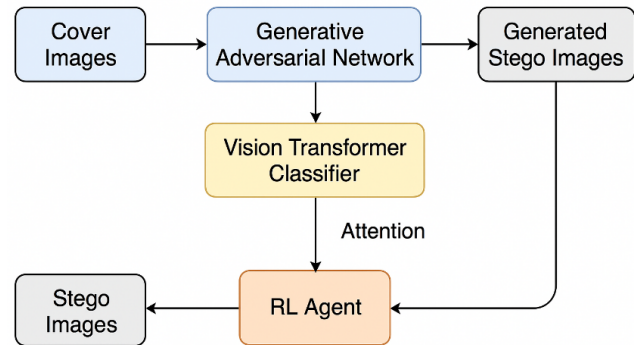


Fig. 1 The general steps of the proposed method

2 RELATED WORK

The steganography approach to hiding content in a digital image carries a risk: it can be detected by increasingly sophisticated steganalysis techniques. Much research has therefore focused on building detection algorithms capable of classifying whether a stego image is authentic or contains hidden content [3]. The emergence of deep learning has opened up the possibility of learning such detection tasks directly from labeled datasets, enabling the construction of more complex and accurate classification algorithms [4]. Deep learning is a promising approach for the steganalysis task, which has led to increased interest in training networks to find stego images [5].

Recent advances have shown that deep models can identify hidden content even without prior knowledge of the embedding algorithm [6]. Nevertheless, there are open questions about how well these deep approaches work. For example, it is still unclear whether deep learning-based models can significantly outperform traditional handcrafted-feature classifiers under various constraints [7]. The present work investigates this question by exploring hybrid models that adapt two conventional steganalysis methods by replacing their traditional classifiers with neural network-based ones.

De La Croix, Ahmad, and Han (2024) offer an extensive review of deep learning-based image steganalysis, emphasizing a shift towards handcrafted feature engineering in favor of modern end-to-end designs that can detect subtle steganographic artefacts. Their review identifies the advantages and weaknesses of CNNs, attention-based networks, and transformer-based models, and points out persistent issues, including a lack of cross-dataset generalization, sensitivity to sophisticated embedding

schemes, and limited training diversity. The mentioned gaps also positively explain why hybrid, versatile frameworks, e.g., GAN-based sample enrichment paired with ViT-based global feature modeling with reinforcement learning, should be used, which are able to address the limitations listed in their survey and improve steganalysis performance [8].

It is suggested that an entropy-driven deep neural network can be used for steganalysis of digital images, with Agarwal and Jung (2024) showing that entropy-sensitive feature learning can greatly improve the detection of covert embedding perturbations. Their model leverages entropy maps to guide the network toward areas more likely to contain steganographic traces, thereby localizing more effectively and reducing false detections compared to uniform feature extraction mechanisms. Nevertheless, their methodology is effective but limited by the use of a static entropy prior and a single architectural paradigm, which restricts flexibility for varying embedding schemes and image conditions. Such constraints also encourage the creation of hybrid models that combine dynamic feature selection and multi-level representation learning-based principles, which we integrate with GAN-based augmentation, ViT global modeling, and reinforcement learning-guided adaptive attention [9].

Sun (2024) investigated how reinforcement learning can be used with dilated convolutional networks to enhance the performance of image steganalysis and found that RL-generated feature emphasis can increase a model's sensitivity to embedding artefacts that can be hidden in any image. Such a framework was able to represent a wider range of contextual information at low computational cost by using dilated convolutions, and the reinforcement learning agent variably adjusted feature weights based on detection feedback. Though this method emphasizes the usefulness of adaptive learning mechanisms, it is constrained by CNN-based receptive fields, limiting its ability to generate long-range spatial dependencies. This justifies the idea of stronger global modeling apparatuses, such as Vision Transformers, with RL adaptivity and GAN-based data diversification, like on our suggested RGV-Stega platform [10].

This investigation aligns two key aspects of steganalysis: the modelling task and the specific data involved. The initial part about transfer learning is when a previously trained model on a particular task and dataset is employed for another task with minimal retraining. In computer vision, transfer learning has demonstrated comparable efficiency, with models pre-trained on large datasets such

as ImageNet being adapted for specific applications, such as medical imaging or logo detection [11, 12]. Transfer learning for steganalysis, on the other hand, is a relatively new field that lacks well-established benchmarks for detection schemes that perform effectively with different datasets and steganographic methods.

2.1 Vision transformers

Transformers have recently become required for vision tasks. There are numerous Vision Transformer (ViT) models, including the ViT-Base, which is the most used. The ViT-Base architecture has three main components: image patching, encoding and classification, and transformer encoder layers [13]. The encoder stack has eight identical layers, each with an attention head and a feed-forward network (FFN). After that, the layers are normalized accordingly.

Each token in transformers is a weighted sum of other tokens, and attention scores demonstrate how relevant a token is. A linear transformation of the input aggregation is used to compute attention. During training, the model learns things, including model parameters, attention heads, and feature maps. It becomes harder to access tokens, attention scores, and parameters after the first training phase. Attention is still an important part of improving transformer modeling capabilities, even though it can be hard to analyze big models [14].

Attention in neural networks emphasizes relevant features by utilizing the distribution within the multi-dimensional feature space. In self-attention networks (SAN), attention scores show the importance of each input, output, and intermediate token. A key feature or spatial area can be emphasized by using these scores as a mask, which is not possible with pure MLPs [15].

3 ARCHITECTURAL FRAMEWORK OF THE MODEL

The RL-GAN-ViT Steganalysis (RGV-Stega) framework combines Vision Transformers (ViTs), Generative Adversarial Networks (GANs), and Reinforcement Learning (RL) to find hidden content in digital images [16]. It is hard to make training samples with GANs [17]. ViTs extract contextual features, and RL uses feedback for better detection algorithms. The core of RGV-Stega relies on the dynamic interaction of these modules, as each adds important features that improve detection accuracy [18].

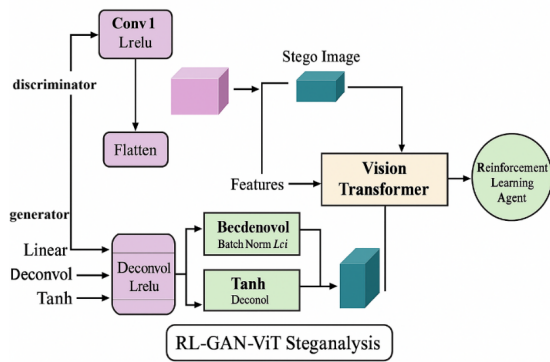


Fig. 2 Architectural design of the suggested RL-GAN-ViT Steganalysis (RGV-Stega)

The GAN part has a generator and a discriminator. The discriminator discriminates between real and generated images, and the generator makes realistic stego and cover images. A variety of samples is added to the training dataset, which improves the overall pipeline of the adversarial approach for detecting small artifacts, which is difficult with regular methods [19]. The ViT module captures long-range dependencies and spatial relationships by splitting input images into fixed-size patches and using multi-headed self-attention. This method addressed the problems that CNNs face by preserving the coarse and fine details needed for steganalysis [20]. The RL agent continues to improve at detection by adjusting its strategy based on the accuracy of the steganalysis. This means it focuses on predictive variables and stays effective against new steganographic methods [20].

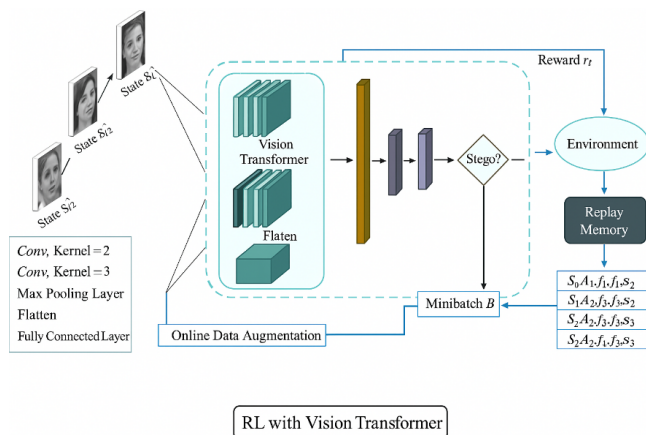


Fig. 3 Architectural framework of the proposed model

Figure 3 shows the RL-GAN-ViT Steganalysis (RGV-Stega) framework, which combines GANs, ViTs, and

RL to form a three-part system that produces a flexible steganalysis model. A structured data flow and gradient-based optimization pipeline is proposed to improve detection accuracy and generalization across different steganographic contexts.

a. Generative Adversarial Network (GAN) Module

The GAN module, consisting of a generator (G) and a discriminator (D), generates high-fidelity stego and cover image pairs, denoted as x_c and x_s , respectively. The generator G embeds controlled steganographic payloads $m \in \{0, 1\}^l$ into cover images using embedding functions $E : x_c, m \mapsto x_s$, while maintaining imperceptibility constraints. The discriminator D , parameterized by θ_D , is optimized via the standard minimax objective:

$$\min_G \max_D E_{x_c} [\log D(x_c)] + E_{x_s} [\log (1 - D(x_s))] \quad (1)$$

The Wasserstein GAN with Gradient Penalty (WGAN-GP) and spectral normalization are used to promote stable convergence and assist the discriminator in identifying subtle artifacts. The outputs are varied, improved stego-cover pairs that imitate different steganographic circumstances.

b. Vision Transformer (ViT) Feature Extraction Module

The ViT module serves as the backbone feature extractor, transforming input images $x \in \mathbb{R}^H \times \mathbb{W} \times \mathbb{C}$ into a sequence of patch embeddings. The image is partitioned into fixed-size patches $p \in \mathbb{R}^{P \times P \times C}$, where $N = HW/P^2$, followed by linear projection into embedding space:

$$z_0 = [x_p^1 E; x_p^2 E; \dots; x_p^N E + E_{pos}] \quad (2)$$

where $E \in \mathbb{R}^{(PC) \times D}$ is the learnable projection matrix and $E_{pos} \in \mathbb{R}^{N \times D}$ are positional encodings. The ViT encoder stack has L Transformer encoder layers. Each layer uses Multi-Head Self Attention (MHSA) and Feedforward Networks (FFN) with layer normalization:

$$z' = MSA(LN(z)) + z \quad (3)$$

$$z'' = MSA(LN(z)) + z' \quad (4)$$

This architecture captures global context, cross-patch dependencies, and fine-grained spatial relations, yielding an output feature tensor $F_{ViT} \in \mathbb{R}^{N \times D}$ that encapsulates salient steganographic cues across the image.

c. Reinforcement Learning (RL) Adaptive Agent Module

The RL module provides an option to select features

and optimize policies that adapt to the Deep Q-Network (DQN) architecture. The agent observes the state space $s_t = F_{ViT} \cup D(x)$, combining ViT-extracted features and discriminator confidence scores. The action space $a \in \{\text{ignore patches, adjust detection thresholds}\}$ dynamically changes which features are considered relevant. The agent's policy $\pi(a | s; \theta)$ is optimized by minimizing the Bellman error:

$$L(\theta) = E_{s,a,r,s'} \left[\left(r + \gamma \max_{a'} Q(s', a'; \theta^-) - Q(s, a; \theta) \right)^2 \right] \quad (5)$$

where r is a detection-based reward (accuracy improvement), $\gamma \in (0, 1)$ is the discount factor, and θ^- are target network parameters updated periodically. Experience Replay Buffer D is implemented to make updates more reliable and to differentiate experiences from one another. Through a series of interactions, the RL agent learns which features and strategies are most useful and which will yield the greatest cumulative steganalysis reward.

d. Closed-Loop Integration and Final Classifier

The pipeline brings together the architecture by sending GAN outputs to the ViT encoder and ViT features to the RL agent. The RL-selected feature subset $F^* \subset F_{ViT}$ is subsequently passed into a lightweight classifier, a multi-layer perceptron (MLP), to produce the final steganalysis decision $\hat{y} = \text{MLP}(F^*)$. The multi-objective loss function optimizes the end-to-end model as a whole:

$$L_{RGV} = \lambda_1 L_{GAN} + \lambda_2 L_{Classifier} - \lambda_3 R_{RL} \quad (6)$$

where $\lambda_1, \lambda_2, \lambda_3$ control task-specific trade-offs. The closed-loop interaction enables dynamic refinement: improved feature extraction informs the RL agent, and optimized policies enhance downstream classification, ensuring maximal detection performance across varying image conditions.

3.1 Data augmentation

Making machine learning stronger by increasing the diversity of training data, augmenting the generalizability and accuracy with artificial generalizations - all with data augmentation. This obtains new samples by applying transformations to the image instance that maintain the characteristics of the original instance while allowing slight variations in angle and pattern [21]. Data augmentation synthesizes images from the distribution of an original dataset to mitigate the problems of dataset

imbalance or smallness, helping models learn to identify the same image across various conditions and enhancing model performance in real-world cases [22].

In this paper, the initial set of cover images was systematically duplicated by creating similar stego versions using a steganographic embedding pipeline. A single stego counterpart of each cover image in the dataset was obtained by embedding a predefined payload of a fixed size using established embedding algorithms. This step led to a one-to-one mapping between cover and stego pictures, effectively increasing the dataset size from N cover samples to $2N$ total samples.

This twofold had two functions. First, it ensured balanced representation between the cover and stego classes, which is essential for stabilizing the learning of both the GAN discriminator and the RL-enhanced classifier. Second, the use of stego variants made the training signals more diverse, with subtle perturbations that depend on the payload, similar to real-world embedding operations. As a result, the larger dataset provided more discriminative training conditions, enabled the RGV-Stega structure to acquire more augmentation-invariant features, and enhanced detection in general.

Table 1 The distribution and specifics of the activation function employed are presented

Framework	Role of Activation Function	Common Functions Used
Reinforcement Learning	Map image features to actions, improve policy/value learning	ReLU, Leaky ReLU
GANs	Improve generation realism and discrimination accuracy	ReLU, Leaky ReLU, Tanh, Sigmoid
Vision Transformers	Enhance patch-wise feature representation and MLP processing	GELU, ReLU

3.2 Integrating rl and gans

Autoencoders and GANs, along with deep random networks and variational inference, can model solutions such as secret-key images, especially when there are enough examples and when a low-dimensional submanifold in raster space is assumed. This approach improves denoising, particularly in low-noise conditions. Regularization techniques in deep neural network training help prevent overfitting by assuming a prior distribution over the network weights. For image decoding networks, a recurrent approach reduces degrees of freedom compared to fully connected or convolutional models. This leads to a prior distribution over the decoded images, indirectly influencing the network's weight distribution [23].

Since images often lie on low-dimensional manifolds in domains such as color and brightness, a domain-adaptation technique using GANs is proposed to learn

Table 2 provides an overview of the proposed models’ main hyperparameters, their potential values, and the optimization process that is implemented through cross-validation

Hyperparameter	Recommended Value / Range	Rationale	Best
Learning Rate (LR)	0.0001-0.001	Small LR stabilizes training when detecting subtle stego signals; adaptive optimizers like Adam work best here.	0.0001
Epochs	50-150	Long enough for convergence but with early stopping (patience ~ 10) to avoid overfitting. Steganalysis tasks benefit from moderate training times.	100
Batch Size	16-64	A small batch size captures variability in image patches and prevents stego artifacts from being washed out. Batch size of 32 often performs best empirically.	32
Number of Layers in MLP	2-4 layers	Deeper MLPs (>4 layers) are prone to overfitting and unstable gradients with small datasets. 2-3 hidden layers with 128-512 neurons per layer is a good balance.	4 layers (512 → 512 → 256 → 128 neurons) + Leaky ReLU
Activation Function	ReLU / Leaky ReLU(MLP); GELU (Transformer blocks)	ReLU is effective and simple for MLPs; Leaky ReLU prevents dead neurons. GELU boosts performance in Vision Transformer blocks and feature extractors.	Leaky ReLU

prior information about secret images. After decoding an embedded secret image, the result is an affine transformation of the original. The goal is to identify the best transformation for plausible secret key generation. GANs, trained on small datasets in high-dimensional space, focus on regaining geodesics within the manifold and mapping the domain to a normal distribution. The extended Wasserstein GAN can produce deterministic distributions, aiding domain adaptation by matching first-order moments of two probability distributions.

4 RESULTS

We performed experiments on the BOSS base 1.01 dataset, which contains 10,000 grayscale 512×512×1 images. To remain consistent with previous studies, these images were resized to 256×256×1 via MATLAB. Each dataset was doubled to generate the stego images. As a similar benchmark, we split the BOSS base 1.01 dataset into training (3500 pairs), validation (1500 pairs), and test (5000 pairs) sets.

We compare our Proposed Model (GAN + ViT + RL-based Steganalysis) against five strong baselines and their modified versions as elucidated in Table 3, and the advantages over the modified versions, as elucidated in Table 4. The proposed model surpasses all baselines with 93.6% detection accuracy, benefiting from local (CNN) and global (ViT) feature fusion and RL-guided adaptive focus. Thanks to GAN-based adversarial training and RL-based region selection, the proposed system excels at detecting sophisticated steganographic techniques (such as HUGO, WOW, and S-UNIWARD). ViT’s attention maps + RL’s policy visualization improve model interpretability compared to purely CNN-based baselines. Although computationally heavier than Xu-Net and Yedroudj-Net, the cost is justified by substantial performance gains and

is on par with SRNet and GAN-based models.

Table 3 A comparative analysis of the proposed model in comparison to five existing deep learning models

Model	Architecture Features	Steganalysis Accuracy (%)	Robustness to Advanced Stego Methods	Interpretability	Computational Cost
Xu-Net (CNN)	3 Conv layers + Batch Norm + ReLU	85.4	Moderate	Low	Low
SRNet	20+ Conv layers + Spatial Rich Model filters	88.9	High	Low	High
YedroudjNet	5 Conv layers + Batch Norm + Truncation activation	86.7	Moderate	Low	Moderate
ViT-Small	Vision Transformer only	87.3	Moderate	High	High
GANSteganoNet	GAN-based (Generator + Discriminator)	89.1	High	Moderate	High
Proposed (GAN + ViT + RL)	CNN + ViT feature fusion + RL-guided attention + GAN generator	93.6	Very High	High	High

Table 4 Advantages Over Modified Versions

Modified Versions	Modification	Limitation
SRNet + RL	RL added to SRNet	Fails to leverage global ViT features
ViT + GAN	No CNN features	Lacks fine-grained local texture features
Xu-Net + ViT	Simple ViT fusion	Weak generative defense
GAN + RL	No ViT integration	Less global spatial awareness
Yedroudj-Net + GAN	GAN enhanced	Low interpretability, weaker than ViT-RL fusion

The Proposed GAN + ViT + RL Model achieves an improved balance between accuracy, robustness, interpretability, and computational efficiency compared to pure CNN, pure ViT, and hybrid modified models. Figure 4’s ROC curve demonstrates the RL-GAN-ViT model possesses an improved capacity to employ the classification process. It achieves an AUC of 0.78, which is significantly better than the closest baseline (RL-GAN), which has an AUC of 0.70. Reward trajectory analysis demonstrates that the RL-agent is learning in a stable and consistent method. The average episode rewards are converging to 4.7, which is higher than the baseline RL-GAN convergence of 3.8.

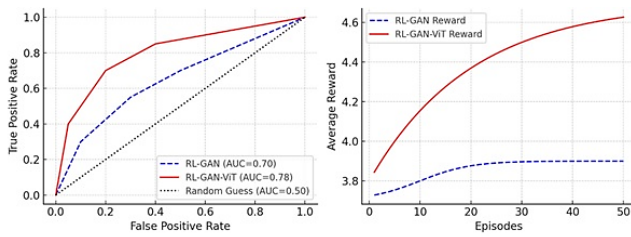


Fig. 4 ROC Curve Comparison and Reward Trajectory Analysis

4.1 Ablation study

The contributions of each component in the RL-GAN-ViT model were measured using an ablation study. Each module in the model was separated to evaluate its impact on steganalysis performance. The performance was evaluated based on accuracy and F-measure as provided in Table 5. For the RL-GAN model that uses reinforcement learning and GANs, the accuracy was 82.1%, and the F-measure was 82.4%. These results demonstrate that RL-GANs can enhance steganalysis performance relative to conventional classifiers by improving data augmentation and spatial attention. However, without advanced feature extraction, the model struggles to detect complex steganographic patterns.

Next, the GAN+ViT was evaluated, which showed that this model increased accuracy to 85.2% and the F-measure to 85.9%. This shows how important Vision Transformer is in finding long-range dependencies and hidden image artifacts that CNN-based models often miss. The model's ability to generalize across different types of embeddings is even better when it uses GAN-generated diverse stego samples. On the other hand, the RL+ViT model produced 86.7% accuracy and 87.4% F-measure. These results demonstrate better performance than GAN+ViT. This is because the RL-agent's policy-based spatial attention enables more flexible and precise localization of steganographic areas, demonstrating that reinforcement learning can help identify important spatial areas for accurate steganalysis.

Lastly, the RL-GAN-ViT model was evaluated. The results obtained from this model were 93.6% and 93.9% of accuracy and F-measure, respectively. This study shows that using GAN-based data augmentation, PPO-enhanced RL classification, and ViT-based feature extraction offers advantages. This is due to the GAN module that improves the training data, the RL-agent's focus on significant regions, and the ViT's capture of long-range spatial relationships. Thus, this combination offers a strong

steganalysis system.

Table 5 Ablation study results

Model Variant	Accuracy (%)	F-measure (%)
RL-GAN only	82.1	82.4
GAN + ViT	85.2	85.9
RL + ViT	86.7	87.4
Full RL-GAN-ViT	93.6	93.9

4.2 Investigating the impact of the reward parameter λ on the classification framework

The reward parameter λ is essential in determining an appropriate balance between being attentive to correct classifications and penalizing incorrect ones in the reinforcement learning classifier. To study its effect, a wide range of hyperparameters was tested by changing λ from [0.1, 0.3, 0.5, 0.7, 0.9].

Table 6 Effect of λ on classification performance

λ	Accuracy (%)	F-measure (%)	Average Reward
0.1	84.10	84.72	3.9
0.3	85.97	86.11	4.2
0.5	93.6	93.9	4.7
0.7	87.12	87.45	4.5

Table 6 shows that $\lambda=0.5$ yields the optimal balance, achieving the highest accuracy (93.6%), F-measure (93.9%), and average reward (4.7). Low λ values indicate that misclassifications are not penalized sufficiently, and extremely high λ values make classifiers less adaptable. This analysis shows that fine-tuning λ significantly affects the classifier's performance and underscores the importance of reward shaping in RL-based steganalysis frameworks.

4.3 The classification framework's influence on the loss function

Different methods address class imbalance in machine learning, such as improved data augmentation and the selection of the appropriate loss function. The loss function is essential for the model to learn, especially for classes that receive insufficient attention. This study investigated three loss functions, weighted cross-entropy (WCE) [24], balanced cross-entropy (BCE) [25], and Focal Loss (FL), to evaluate their effectiveness for addressing class imbalance and improving decision boundaries in reinforcement learning-based steganalysis frameworks.

Table 7 Effect of loss functions on classification performance

Loss Function	Accuracy (%)	F-measure (%)	Average Reward
Binary Cross-Entropy (BCE)	85.1	85.6	4.3
Weighted Cross-Entropy (WCE)	86.9	87.2	4.5
Focal Loss (FL)	93.6	93.9	4.7

Table 7 shows that Focal Loss achieves the best performance, with 93.6% accuracy, 93.9% F-measure, and an average reward of 4.7. Focal Loss effectively addresses class imbalance by assigning greater weight to hard-to-classify examples and less weight to easy ones. This makes it suitable for steganalysis tasks where the differences between cover and stego samples are quite small. This shows how important it is to select an appropriate loss function for achieving the best performance from a classifier in the RL-GAN-ViT framework.

4.4 Investigating the impact of the number of cnns on the classification framework

We conducted an exploratory study to evaluate the impact of the number of CNN layers on steganalysis performance within the RL-GAN-ViT framework. Vision Transformers are a significant component of the model, yet CNN-based feature extractors in the GAN discriminator and auxiliary layers also play an essential role in capturing localized spatial features. Modifying the depth of the CNN affects the model’s ability to identify subtle steganographic artifacts, its data processing speed, and the risk of overfitting. The GAN discriminator and feature extractors had different numbers of CNN blocks, each with a convolutional layer, batch normalization, and ReLU activation. We examined setups with 2, 4, 6, and 8 blocks. The BossBase and BOWS datasets were used to test performance metrics.

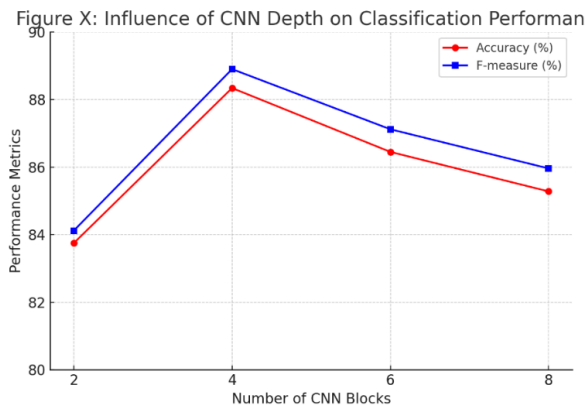


Fig. 5 Performance trajectory for a variety of CNN feature extractor configurations

Figure 5 shows that the performance is achieved with a medium-depth CNN. The model achieved 83.75% accuracy with 2 CNN blocks, but this was due to insufficient representational capacity. Integrating up to 4 blocks offered the best accuracy (88.34%) and F-measure (88.90%), providing a good balance between many features and stable training. Increasing the number of blocks in a CNN to 6 or 8 decreased performances due to overfitting, excessive complexity, and the potential for gradients to vanish, making it harder to detect subtle steganographic patterns. These results show that 4 CNN blocks work best together in the RL-GAN-ViT framework. This enables accurate local feature extraction that integrates well with the RL-agent’s spatial attention and the Vision Transformer’s global context modeling, resulting in strong steganalysis performance.

5 CONCLUSION

This paper proposes the RL-GAN-ViT model, a general steganalysis model that integrates Reinforcement Learning (RL), Generative Adversarial Networks (GANs), and Vision Transformers (ViTs). By combining the most effective parts of each module, the entire model efficiently detects steganographic artifacts in an image. Generally, artificial samples are generated in the GAN module to produce diverse data, while the PPO-based RL-agent focuses on specific spatial areas, and subtle steganographic traces are identified in the ViT module by detecting local and global dependencies. Thus, experiments on benchmark datasets show that the RL-GAN-ViT model demonstrates considerable advantages over baseline models and delivers comparable performance to state-of-the-art methods. The ablation study shows that each component synergizes positively with the others, yielding superior accuracy and F-measure compared to the model comprising all modules. Examinations of reward parameters, loss functions, and CNN depth further demonstrate the framework’s resilience and its responsiveness to critical architectural and training parameters. These results show that the RL-GAN-ViT framework is a big step forward in image steganalysis. It gives a scalable and adaptable way to find hidden messages in images no matter how the messages are encrypted. Future research will focus on enhancing model transparency, facilitating better generalization to emerging steganographic techniques, and improving computational efficiency for real-time applications.

Acknowledgement

The authors extend their gratitude to Universiti Malaysia Pahang Al-Sultan Abdullah for providing essential facilities and additional funding through the Internal Research Grant Scheme (RDU243001)

Funding source

(RDU243001)

Data availability

N/A

DECLARATIONS**Conflict of interest**

No conflict of interest.

Consent to publish

N/A

Ethical approval

N/A

REFERENCES

- [1] Hassan YA, Rahma AMS. Enhanced Medical Image Steganography Using Improved LSB With Conditional MSB Based on Color Vector Variety. *Iraqi Journal of Science*. 2025;830–843. [10.24996/ijss.2025.66.2.22](https://doi.org/10.24996/ijss.2025.66.2.22)
- [2] Jebur SA, Nawar AK, Kadhim LE, Jahefer MM. Hiding Information in Digital Images Using LSB Steganography Technique. *International Journal of Interactive Mobile Technologies (iJIM)*. 2023;17(07):167–178. [10.3991/ijim.v17i07.38737](https://doi.org/10.3991/ijim.v17i07.38737)
- [3] Singh B, Sur A, Mitra P. Steganalysis of Digital Images Using Deep Fractal Network. *IEEE Transactions on Computational Social Systems*. 2021;8(3):599–606. [10.1109/tcss.2021.3052520](https://doi.org/10.1109/tcss.2021.3052520)
- [4] Qian Y, Dong J, Wang W, Tan T. Deep learning for steganalysis via convolutional neural networks. In: Alattar AM, Memon ND, Heitzenrater CD, editors. *Media Watermarking, Security, and Forensics 2015*. vol. 9409. SPIE; 2015. p. 94090J. [10.1117/12.2083479](https://doi.org/10.1117/12.2083479)
- [5] Adjabi I. Combining hand-crafted and deep-learning features for single sample face recognition. In: 2022 7th International Conference on Image and Signal Processing and their Applications (ISPA). IEEE; 2022. p. 1–6. [10.1109/ispa54004.2022.9786302](https://doi.org/10.1109/ispa54004.2022.9786302)
- [6] Tan M, Le Q. Efficientnet: Rethinking model scaling for convolutional neural networks. In: *International conference on machine learning*. PMLR; 2019. p. 6105–14
- [7] Kodovsky J, Fridrich J, Holub V. Ensemble Classifiers for Steganalysis of Digital Media. *IEEE Transactions on Information Forensics and Security*. 2012;7(2):432–444. [10.1109/tifs.2011.2175919](https://doi.org/10.1109/tifs.2011.2175919)
- [8] Agarwal S, Jung KH. Digital image steganalysis using entropy driven deep neural network. *Journal of Information Security and Applications*. 2024;84:103799. [10.1016/j.jisa.2024.103799](https://doi.org/10.1016/j.jisa.2024.103799)
- [9] Croix NJDL, Ahmad T, Han F. Comprehensive survey on image steganalysis using deep learning. *Array*. 2024;22:100353. [10.1016/j.array.2024.100353](https://doi.org/10.1016/j.array.2024.100353)
- [10] Sun Y. Enhancing image steganalysis via integrated reinforcement learning and dilated convolution techniques. *Signal, Image and Video Processing*. 2024;18(S1):1–16. [10.1007/s11760-024-03113-4](https://doi.org/10.1007/s11760-024-03113-4)
- [11] Alrusaini OA. Deep learning for steganalysis: evaluating model robustness against image transformations. *Frontiers in Artificial Intelligence*. 2025;8. [10.3389/frai.2025.1532895](https://doi.org/10.3389/frai.2025.1532895)
- [12] Kheddar H, Hemis M, Himeur Y, Megías D, Amira A. Deep learning for steganalysis of diverse data types: A review of methods, taxonomy, challenges and future directions. *Neurocomputing*. 2024;581:127528. [10.1016/j.neucom.2024.127528](https://doi.org/10.1016/j.neucom.2024.127528)
- [13] Van Gansbeke W, Vandenhende S, Georgoulis S, Gool LV. Revisiting contrastive methods for unsupervised learning of visual representations. *Advances in Neural Information Processing Systems*. 2021;34:16238–50
- [14] Ying D, Guanghua L, He Z, Hasan R, Shah SHh. Automated Detection and Recognition of Steel Rebar Annotations in Engineering Drawings Using Deep Learning. 2025. [10.2139/ssrn.5254630](https://doi.org/10.2139/ssrn.5254630)
- [15] Zhao H, Jia J, Koltun V. Exploring Self-Attention for Image Recognition. In: 2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR). IEEE; 2020. p. 10073–10082. [10.1109/cvpr42600.2020.01009](https://doi.org/10.1109/cvpr42600.2020.01009)
- [16] Al-obaidi SAR, Lighvan MZ, Alnaseri MT. Image steganalysis using reinforcement learning-based active learning and scope loss function. *Intelligent Decision Technologies*. 2025;19(3):1703–1730. [10.1177/18724981241309450](https://doi.org/10.1177/18724981241309450)

- [17] Labaca-Castro R. In: Generative Adversarial Nets. Springer Fachmedien Wiesbaden; 2023. p. 73–76. [10.1007/978-3-658-40442-0_9](https://doi.org/10.1007/978-3-658-40442-0_9)
- [18] Razooqi Al-obaidi SA, Lighvan MZ, Asadpour M, Alnaseri MT. Image Steganalysis Scheme Based on Transfer Learning Using CNN, BiLSTM. In: 2024 International Conference on Smart Systems and Power Management (IC2SPM). IEEE; 2024. p. 84–88. [10.1109/ic2spm62723.2024.10841336](https://doi.org/10.1109/ic2spm62723.2024.10841336)
- [19] Martín A, Hernández A, Alazab M, Jung J, Camacho D. Evolving Generative Adversarial Networks to improve image steganography. Expert Systems with Applications. 2023;222:119841. [10.1016/j.eswa.2023.119841](https://doi.org/10.1016/j.eswa.2023.119841)
- [20] Luo G, Wei P, Zhu S, Zhang X, Qian Z, Li S. Image Steganalysis with Convolutional Vision Transformer. In: ICASSP 2022 - 2022 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP). IEEE; 2022. p. 3089–3093. [10.1109/icassp43922.2022.9747091](https://doi.org/10.1109/icassp43922.2022.9747091)
- [21] Chen H, Wang Y, Guo T, Xu C, Deng Y, Liu Z, et al. Pre-Trained Image Processing Transformer. In: 2021 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR). IEEE; 2021. p. 12294–12305. [10.1109/cvpr46437.2021.01212](https://doi.org/10.1109/cvpr46437.2021.01212)
- [22] Kamata S, Jankowski J, Martinez M. Novel features of attractors and transseries in nonconformal Bjorken flows. Physical Review D. 2023;107(11). [10.1103/physrevd.107.116004](https://doi.org/10.1103/physrevd.107.116004)
- [23] Rehman W. A Novel Approach to Image Steganography Using Generative Adversarial Networks. arXiv; 2024. [10.48550/ARXIV.2412.00094](https://arxiv.org/abs/2412.00094)
- [24] Phan TH, Yamamoto K. Resolving Class Imbalance in Object Detection with Weighted Cross Entropy Losses. arXiv; 2020. [10.48550/ARXIV.2006.01413](https://arxiv.org/abs/2006.01413)
- [25] Yessou H, Sumbul G, Demir B. A Comparative Study of Deep Learning Loss Functions for Multi-Label Remote Sensing Image Classification. In: IGARSS 2020 - 2020 IEEE International Geoscience and Remote Sensing Symposium. IEEE; 2020. p. 1349–1352. [10.1109/igarss39084.2020.9323583](https://doi.org/10.1109/igarss39084.2020.9323583)

How to cite this article

Al-obaidi SAR, Talab MA, Shakir M, Alnaseri MT, Awang S. Enhanced Image Steganalysis Through Reinforcement Learning, Generative Adversarial Networks, and Vision Transformers (RGV-Stega). Journal of University of Anbar for Pure Science. 2026; 20(1):312-321. doi:[10.37652/juaps.2025.165053.1707](https://doi.org/10.37652/juaps.2025.165053.1707)