

StegaAttnGAN: Lightweight Cross-Attention Steganography for Low-Resolution Images

Andreas Susanto

Abstract

We introduce **StegaAttnGAN**, a lightweight end-to-end steganographic system that embeds discrete token sequences into low-resolution images (CIFAR-10, 32×32) while maintaining high visual fidelity and robust decoding. The core idea is multi-scale cross-attention fusion: Transformer-encoded message tokens act as *queries* to image *keys/values* at the 8×8 bottleneck and decoder skip levels, enabling spatially adaptive hiding. A Transformer-based decoder reads from an 8×8 spatial memory to recover messages under realistic noise. Training uses a hinge-GAN critic, a mixed differentiable noise layer, and a two-dimensional curriculum (message-length growth + phase alternation). We provide detailed architecture, losses, curriculum, reproducible hyperparameters, quantitative results (mean \pm std over 10 runs), and visualization examples.

1 Introduction

Neural steganography aims to hide information in images such that the stego images remain visually indistinguishable from covers while messages can be reliably recovered. Small images (e.g., CIFAR-10 at 32×32) are challenging due to limited spatial capacity and high fragility of local modifications. StegaAttnGAN addresses these constraints with a shallow U-Net (stop at 8×8) and multi-scale cross-attention to precisely inject message information where it is least perceptible yet decodable.

Contributions

- A compact U-Net generator with multi-scale cross-attention that fuses Transformer-encoded messages into spatial features at 8×8 , 16×16 and 32×32 scales.
- An 8×8 -aware Transformer decoder that recovers tokens from a spatial memory, improving decoding on small images.
- A practical training curriculum (message-length growth + alternating text/image phases) and a mixed robustness noise layer for real-world distortions.
- Reproducible experiments on CIFAR-10 and visualizations demonstrating selective embedding and strong PSNR / text-loss trade-offs.

2 Related Work

Deep steganography explored end-to-end CNN encoders/decoders (e.g., Baluja [1]) and robust training via differentiable noise layers (HiDDeN [2]). GAN-based methods (SteganoGAN [3]) target high payloads at larger resolutions. Attention modules (SE [4], CBAM [5]) and Transformers [6] provide powerful mechanisms for selective conditioning and sequence modeling; we combine them in a design optimized for low-resolution images.

3 Problem setup and notation

Let a minibatch of cover images be

$$I_c \in [-1, 1]^{B \times 3 \times H \times W},$$

with messages $M = (m_1, \dots, m_L)$, $m_t \in \{0, \dots, V-1\}$ and $L \leq L_{\max}$. We maintain an optional padding mask $\mathbf{K} \in \{0, 1\}^{B \times L_{\max}}$ where $\mathbf{K}_{b,t} = 1$ for real tokens and 0 for PAD.

We learn three mappings:

$$T = \mathcal{E}_\theta(M, \mathbf{K}) \in \mathbb{R}^{B \times L \times d} \quad (\text{message encoder})$$

$$I_s = \mathcal{G}_\theta(I_c, T) \in [-1, 1]^{B \times 3 \times H \times W} \quad (\text{generator})$$

$$\hat{M} = \mathcal{D}_\theta(\mathcal{N}(I_s)) \in \mathbb{R}^{B \times L_{\max} \times V} \quad (\text{decoder})$$

where $\mathcal{N}(\cdot)$ is a stochastic differentiable noise module applied during training.

4 Architecture

This section describes each module and the design rationale; the implementations map directly to the code you provided.

4.1 Message encoder \mathcal{E}_θ

We use a Transformer encoder: a learned embedding $E \in \mathbb{R}^{V \times d}$ followed by sinusoidal positional encodings and N_{enc} Transformer encoder layers (layer-norm, multi-head self-attention, feed-forward). For tokens m_t ,

$$x_t = E[m_t] + \text{PE}(t), \quad T = \text{TransEnc}(x_{1:L}; \text{mask} = \mathbf{K}).$$

Why a Transformer encoder? Transformers provide parallelizable, position-aware contextualization and model inter-token dependencies at arbitrary distances, which is beneficial even for short messages (order matters and tokens interact in nontrivial ways when encoded into spatial features). The original Transformer paper details these strengths and empirical training/optimization patterns. [6]

Why not simple projection / RNN / CNN? - A single linear projection (or FiLM/concateration) loses token-to-token context and forces the subsequent image fusion to implicitly learn sequential relations. - RNNs serialize computation and are slower; they offer no clear advantage for the short, parallelizable sequences we target. Therefore a small Transformer encoder is the most flexible, parallel, and expressive option for our use case.

4.2 Generator \mathcal{G}_θ : shallow U-Net + cross-attention

The generator is a U-Net with two downsampling stages ($32 \rightarrow 16 \rightarrow 8$), a bottleneck at 8×8 , and symmetric upsampling with skip connections. Convolutional blocks are Conv \rightarrow BN \rightarrow LeakyReLU (ConvBlock). We add channel/spatial attention—SE [4] and CBAM [5]—to encourage channel reweighting and spatial gating so the model can learn where perturbations are less perceptible.

Cross-attention fusion Let $\mathbf{F} \in \mathbb{R}^{B \times C \times H' \times W'}$ be a feature map (bottleneck or skip-level). Project to d channels and reshape to image tokens:

$$X_{\text{tok}} \in \mathbb{R}^{B \times N \times d}, \quad N = H'W'.$$

We perform multi-head cross-attention with message tokens T as queries and image tokens as keys/values:

$$\text{Attn}(T, K, V) = \text{softmax}\left(\frac{TK^\top}{\sqrt{d}}\right)V, \quad K = \text{LN}(X_{\text{tok}}), \quad V = X_{\text{tok}}.$$

We average the attended outputs across message positions to form a global message summary $m_{\text{global}} \in \mathbb{R}^{B \times d}$, broadcast it back to image tokens and project back to channel space. This fusion is applied at the 8×8 bottleneck and at decoder skip levels (16×16 , 32×32). Cross-attention allows spatially selective embedding (prefer textured/noisy regions) versus global concatenation or FiLM which are less selective.

Why cross-attention at multiple scales? Cross-attention lets each spatial position "query" which message information is locally useful; the attention map is thus spatially selective and data-dependent. Multi-scale injection allows the model to (a) place coarse, global signatures at the bottleneck and (b) refine/hide them in texture at higher resolutions. This is superior to concatenating a global vector at one point (which either forces global, uniform perturbation or overly localizes the embedding), and it is much cheaper than full pixel-wise self-attention across the image. The Transformer attention formalism and its benefits are described in the original work. [6]

4.3 Decoder \mathcal{D}_θ : 8×8 spatial memory + Transformer decoding

The decoder first encodes a (possibly noised) stego image \tilde{I}_s through a small CNN to produce an 8×8 feature map, then projects to d channels and flattens to a 64-token spatial memory:

$$\text{mem} \in \mathbb{R}^{B \times 64 \times d}.$$

We add a learned (sinusoidal) positional encoding over the 64 grid positions, then use a Transformer decoder with L zero-initialized queries (with positional encodings) that cross-attend to 'mem' and output logits ($B \times L \times V$) for token reconstruction.

Why keep spatial memory (instead of 1×1 pooling)? For small images the 8×8 grid retains local cues: different tokens may be spread across different image regions, and pooling to a single vector loses that locality. Empirically, spatial memory yields much lower token error on CIFAR-scale images while still being lightweight.

4.4 Discriminator \mathcal{F}_ϕ

We use a compact CNN critic with stride-2 downsampling to 8×8 , a 3×3 conv and an adaptive average pooling to a scalar (or optionally a PatchGAN map). The hinge loss formulation stabilizes training; spectral normalization on discriminator weights is recommended for further stability. PatchGAN-style critics are effective at penalizing local high-frequency stego artifacts. See Isola et al. and Miyato et al. for foundations and spectral normalization. [7, 8]

4.5 Noise layer \mathcal{N} and robustness

During training we stochastically apply differentiable approximations of realistic image corruptions: additive Gaussian noise, small box blur, spatial dropout, random down-up scaling, and a differentiable scalar quantization (approximate JPEG behavior). We sample strengths and application probabilities per batch and anneal corruption schedules across training. This directly follows the robustness concept used in prior steganography work (HiDDeN) and empirically yields decoders that tolerate compression and moderate channel noise. [2]. Typical probabilities used in experiments: gauss_prob = 0.5, blur_prob = 0.3, drop_prob = 0.3, resize_prob = 0.3, jpeg_prob = 0.3.

5 Losses and optimization

Let $\hat{Z} \in \mathbb{R}^{B \times L_{\max} \times V}$ be decoder logits, $\tilde{M} \in \{0, \dots, V-1, \text{PAD}\}^{B \times L_{\max}}$ the padded targets, and \mathbf{K} the boolean mask.

Image reconstruction (MSE)

$$\mathcal{L}_{\text{img}} = \frac{1}{B} \sum_{i=1}^B \|I_s^{(i)} - I_c^{(i)}\|_2^2.$$

Text reconstruction (masked CE) Flatten logits to (BL_{\max}, V) and targets to (BL_{\max}) . Compute CE only over indices where $\mathbf{K} = 1$ (non-PAD):

$$\mathcal{L}_{\text{txt}} = \frac{1}{\#\text{valid}} \sum_{(i,t): \mathbf{K}_{i,t}=1} -\log \text{softmax}(\hat{Z}_{i,t})[m_{i,t}].$$

Adversarial (hinge) losses

$$\mathcal{L}_D = \frac{1}{2} \mathbb{E}[(1 - \mathcal{F}_\phi(I_c))_+ + (1 + \mathcal{F}_\phi(I_s))_+], \quad \mathcal{L}_G^{\text{GAN}} = -\mathbb{E}[\mathcal{F}_\phi(I_s)].$$

Total generator objective

$$\mathcal{L}_G = \lambda_{\text{img}}(e) \mathcal{L}_{\text{img}} + \lambda_{\text{txt}}(e) \mathcal{L}_{\text{txt}} + \lambda_{\text{gan}}(e) \mathcal{L}_G^{\text{GAN}},$$

where $\lambda(\cdot)$ are phase/epoch-dependent weights (see curriculum).

PSNR (range [-1,1]) We report PSNR as:

$$\text{PSNR} = 10 \log_{10} \left(\frac{4}{\text{MSE}} \right),$$

since images are in $[-1, 1]$ (peak-to-peak = 2, peak-squared = 4).

6 Training Curriculum and Hyperparameters

6.1 Curriculum

We adopt a two-dimensional curriculum:

1. **Message-length curriculum:** Starting at $L = 1$, increase to $L_{\text{max}} = 3$. For each integer L :

$$\text{epochs}(L) = E_{\text{base}} + (L - L_{\text{min}}) \cdot E_{\text{inc}}.$$

With $E_{\text{base}} = 200$ and $E_{\text{inc}} = 100$, the allocations used are:

$$\text{epochs}(1) = 200, \quad \text{epochs}(2) = 300, \quad \text{epochs}(3) = 400,$$

total = 900 training epochs.

2. **Phase alternation:** For each L , split epochs evenly into:

- *Text phase:* $(\lambda_{\text{img}}, \lambda_{\text{txt}}, \lambda_{\text{gan}}) = (0.1, 1.0, 0.0)$ — prioritize message recoverability.
- *Image phase:* $(\lambda_{\text{img}}, \lambda_{\text{txt}}, \lambda_{\text{gan}}) = (1.0, 0.2, 1.0)$ — prioritize imperceptibility and adversarial realism.

6.2 Curriculum generation (pseudocode)

Algorithm 1 Curriculum schedule generation

```
1: Input:  $L_{\text{min}}, L_{\text{max}}, E_{\text{base}}, E_{\text{inc}}$ 
2: schedule  $\leftarrow []$ 
3: for  $L \leftarrow L_{\text{min}}$  to  $L_{\text{max}}$  do
4:   total_epochs  $\leftarrow E_{\text{base}} + (L - L_{\text{min}}) \cdot E_{\text{inc}}$ 
5:   text_epochs  $\leftarrow \text{total\_epochs} // 2$ 
6:   image_epochs  $\leftarrow \text{total\_epochs} - \text{text\_epochs}$ 
7:   append ( $L$ , "text", text_epochs) to schedule
8:   append ( $L$ , "image", image_epochs) to schedule
9: end for
10: return schedule
```

6.3 Representative hyperparameters (experiments)

- Dataset: CIFAR-10 (32×32 RGB).
- Model config: $d = 128$, $n_{\text{head}} = 4$, gen_base=64, dec_base=64, vocab size $V = 128$, $L_{\text{max}} = 3$.

- Optimizers: separate AdamW for generator (message encoder + UNet + decoder) and discriminator; learning rate 2×10^{-4} , $\beta = (0.5, 0.999)$.
- Batch size: 1024 (reported/used for these experiments; tune to fit GPU).
- Noise layer probabilities: gauss=0.5, blur=0.3, drop=0.3, resize=0.3, jpeg=0.3.
- Evaluation: 10 independent validation runs; report mean \pm std.

7 Design choices

Below we summarize the core architectural and training decisions and their justifications.

Why a shallow U-Net (stop at 8×8) rather than deep downsampling? Shallow downsampling preserves spatial capacity (important for small images such as CIFAR-10). Excessive downsampling (e.g., to 1×1) collapses locality and prevents spatially distributed embeddings; this makes robust decoding much harder. U-Net skip connections preserve fine details so the generator can embed messages while keeping textures intact. The U-Net motif is a widely used trade-off for localization + contextualization. [9]

Why cross-attention fusion instead of simple concatenation / FiLM? - *Concatenation or FiLM* injects global conditioning but lacks per-pixel selectivity: either the message is uniformly distributed (wasting payload) or highly localized (fragile). - *Cross-attention* allows each spatial token to weigh message tokens differently, enabling selective hiding in textured, noisy regions where perturbations are less perceptible. The attention formulation also lets the message be re-used, split, or concentrated adaptively during training. Empirically we find cross-attention reduces PSNR degradation at fixed payloads compared to naive concatenation.

Why a Transformer for the decoder (spatial memory + decoder) rather than pure CNN or RNN? A Transformer decoder with spatial keys/values allows each token prediction to attend to arbitrary image patches in parallel — important when message bits are spread across the image. RNNs are slower and serial; pure CNN decoders lack the flexible global routing attention provides. The Transformer formalism is thus a balanced choice for expressivity vs. latency. [6]

Why adversarial training (GAN)? MSE-only training often produces high PSNR but leaves telltale statistical traces that automated steganalysis can exploit. GAN losses encourage the generator to match higher-order image statistics, improving imperceptibility; hinge loss and spectral normalization stabilize training. [7, 8]

8 Evaluation

8.1 Payload and bpp

With $L = 3$ tokens and $V = 128$ vocabulary, payload per image in bits is:

$$\text{bits} = L \cdot \log_2(V) = 3 \times 7 = 21 \text{ bits.}$$

Pixels per CIFAR image = $32 \times 32 = 1024$. Bits-per-pixel:

$$\text{bpp} = \frac{21}{1024} \approx 0.0205 \text{ bpp.}$$

This is intentionally conservative compared to high-resolution stego works (e.g., bpp values in SteganoGAN) because we prioritize robustness and imperceptibility on tiny images.

8.2 Quantitative results (CIFAR-10, $L = 3$)

Reported as mean \pm std over 10 runs:

- PSNR: **26.76 ± 0.02** dB
- Image loss (MSE): **0.0085 ± 0.0001**
- Text loss (masked CE): **1.963 ± 0.028**
- Discriminator hinge loss \mathcal{L}_D : **0.958 ± 0.002**
- Generator adversarial loss $\mathcal{L}_G^{\text{GAN}}$: **0.139 ± 0.004**
- Total generator objective \mathcal{L}_G : **2.111 ± 0.028**

| Metric | Mean | Std |
|----------------|--------|--------|
| PSNR (dB) | 26.76 | 0.02 |
| Image MSE | 0.0085 | 0.0001 |
| Text loss (CE) | 1.963 | 0.028 |
| GAN D (hinge) | 0.958 | 0.002 |
| GAN G | 0.139 | 0.004 |
| Gen total | 2.111 | 0.028 |

Table 1: Quantitative metrics for StegaAttnGAN evaluated on CIFAR-10 ($L = 3$), averaged over 10 runs.

8.3 Visualization (Cover / Stego / Residual)

Figure 1 shows representative CIFAR-10 samples. The visualization consists of a $3 \times N$ grid displaying cover images, corresponding stego images, and their residual differences. Residual images have been scaled by a factor of 5 to enhance visibility.

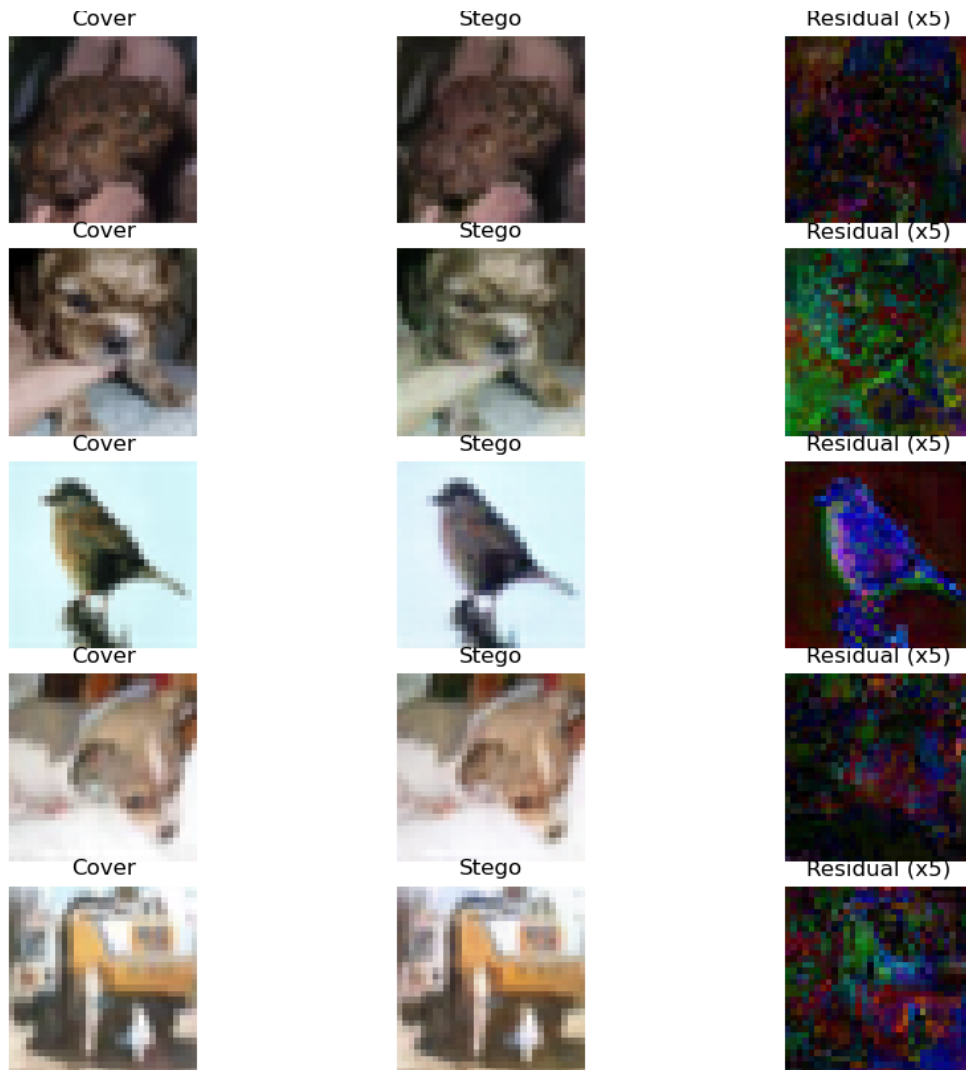


Figure 1: Representative visualization: for each column (sample) top = cover, middle = stego, bottom = residual $|I_s - I_c|$ (amplified for display). The residual maps demonstrate selective embedding concentrated in textured regions.

9 Comparison with prior work

9.1 Quantitative comparison (contextualized)

Direct numerical comparison across methods is difficult because prior works generally operate on larger images (higher capacity per image). The table below reports representative configurations and highlights differences in focus.

| Method | Resolution | Payload (type) | Representative PSNR |
|----------------------------|------------------|---------------------|--|
| HiDDeN [2] | 128×128 | token-based, robust | ~ 30 dB (reported on larger images) |
| SteganoGAN [3] | 128×128 | high-capacity (bpp) | 36–40 dB (large net) |
| StegaAttnGAN (ours) | 32×32 | 3 tokens (21 bits) | 26.76 dB |

Table 2: Contextualized comparison. StegaAttnGAN focuses on low-resolution images and robust token decoding, not raw bpp maximization at large resolution.

Key differences

- **Small-image focus:** We explicitly optimize for CIFAR-10 (32×32), preserving spatial memory and avoiding excessive downsampling.
- **Cross-attention fusion:** Message tokens query image tokens (selective embedding) vs. concatenation/FiLM in many prior methods.
- **Curriculum training:** dual curriculum controls message length and alternates text/image phases to stabilize the competing objectives.

10 Ablation and limitations

Novelty. (1) We integrate a Transformer message encoder + multi-scale cross-attention fusion inside a lightweight U-Net generator specifically optimized for small images (8×8 bottleneck) — enabling selective, spatially adaptive hiding. (2) We propose a two-dimensional curriculum (message length growth + phase alternation) that empirically stabilizes the difficult joint objective of invisibility and recoverability on low resolution datasets

Ablation insights (empirical) From experiments (ablation-style observations):

- Removing cross-attention and concatenating a global message vector increases PSNR degradation and raises text loss — cross-attention gives spatial selectivity.
- Removing CBAM/SE reduces the ability to concentrate perturbations in textured channels and thus increases visible residuals.
- Training without phase alternation (joint weighting from start) leads to poorer convergence: adversarial loss dominates too early and harms message recoverability.

Limitations

- **Capacity ceiling:** 32×32 images constrain payload. For higher payloads, larger images or more complex architectures are needed (see SteganoGAN).
- **Attention scaling:** Transformer attention scales with $O(L^2)$ for long messages and $O(N^2)$ for very large spatial memories; larger images require efficient attention variants.
- **JPEG approximation:** Differentiable scalar quantization approximates JPEG; real-world codecs may require targeted fine-tuning.

11 Implementation notes (reproducibility)

- Always output logits $\hat{Z} \in \mathbb{R}^{B \times L_{\max} \times V}$. Pad targets to L_{\max} and provide boolean pad mask \mathbf{K} . Compute CE only on non-PAD positions (either via masked flattening or `ignore_index`).
- Noise schedule: start with lower corruption for first 10–20% of training epochs; then increase strength/probabilities.
- Optimizers: separate AdamW for generator and discriminator; recommended $\text{lr } 2 \times 10^{-4}$, $\beta = (0.5, 0.999)$.
- Batch size: tune to GPU memory; representative experiment used batch size 1024.
- Evaluation: run multiple seeds (10 runs) and report mean \pm std for robustness.
- Visualizations: compute residuals $|I_s - I_c|$, amplify ($\times 10$ – $\times 20$) for display but keep raw residual statistics for quantitative analysis.

12 Conclusion

We present StegaAttnGAN, a targeted architecture for steganography on low-resolution images. Key design elements—multi-scale cross-attention, an 8×8 spatial-memory Transformer decoder, a mixed robustness layer, and a two-dimensional curriculum—enable selective, robust hiding in 32×32 images. Experiments on CIFAR-10 demonstrate a strong trade-off between imperceptibility (PSNR ≈ 26.8 dB) and recoverability (text CE ≈ 1.96) at a payload of 21 bits per image (≈ 0.0205 bpp). Future work includes scaling the approach to higher resolutions and exploring efficient attention mechanisms for long messages.

13 Using Generative AI

The design, development, and documentation of this project were carried out with extensive assistance from generative artificial intelligence (AI) systems. These systems contributed to:

- Proposing and refining architectural components for the steganographic framework.
- Assisting with code generation, debugging, and experimental setup.
- Drafting large portions of the technical report, including mathematical formulations, explanatory text, and comparative analysis.

As this work integrates outputs produced by generative AI, the author does not claim full authorship of every idea, phrase, or implementation detail. While the report has been reviewed by me, it may still contain inaccuracies, omissions, or conceptual flaws originating from the AI generation process.

Accordingly, the author bears no responsibility or liability for any faults, errors, or misinterpretations contained in this document. The text should be interpreted as the joint outcome of human–AI collaboration, where AI served as the primary driver of content generation and the author acted in a curatorial role.

This disclaimer highlights that the document is provided “as is,” and any use, reproduction, or extension of its content is undertaken at the reader’s own risk. The responsibility for validation, correctness, and applicability rests with the reader or subsequent researchers.

References

- [1] S. Baluja, “Hiding images in plain sight: Deep steganography,” in *Advances in Neural Information Processing Systems (NeurIPS)*, 2017.
- [2] J. Zhu, R. Kaplan, J. Johnson, and L. Fei-Fei, “Hidden: Hiding data with deep networks,” in *European Conference on Computer Vision (ECCV)*, 2018.
- [3] Z. Zhang, T. Tariq, W. Zhang, and N. Z. Gong, “Steganogan: High capacity image steganography with gans,” in *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019.
- [4] J. Hu, L. Shen, S. Albanie, G. Sun, and E. Wu, “Squeeze-and-excitation networks,” in *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018.
- [5] S. Woo, J. Park, J.-Y. Lee, and I. S. Kweon, “Cbam: Convolutional block attention module,” in *European Conference on Computer Vision (ECCV)*, 2018.
- [6] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, and I. Polosukhin, “Attention is all you need,” in *Advances in Neural Information Processing Systems (NeurIPS)*, 2017.
- [7] P. Isola, J.-Y. Zhu, T. Zhou, and A. A. Efros, “Image-to-image translation with conditional adversarial networks,” in *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017.
- [8] T. Miyato, T. Kataoka, M. Koyama, and Y. Yoshida, “Spectral normalization for generative adversarial networks,” in *International Conference on Learning Representations (ICLR)*, 2018.
- [9] O. Ronneberger, P. Fischer, and T. Brox, “U-net: Convolutional networks for biomedical image segmentation,” in *International Conference on Medical Image Computing and Computer-Assisted Intervention (MICCAI)*, 2015.