




Lightweight Deep Neural Network–Based Image Style Transfer for Pattern Design

Wentao Huo 

Department of Visual Communication Design, Hebei Art & Design Academy, Baoding 071000, China

Corresponding Author Email: 15930926511@163.com

Copyright: ©2026 The author. This article is published by IIETA and is licensed under the CC BY 4.0 license (<http://creativecommons.org/licenses/by/4.0/>).

<https://doi.org/10.18280/ts.430221>

ABSTRACT

Received: 7 November 2025

Revised: 30 January 2026

Accepted: 19 February 2026

Available online: 30 April 2026

Keywords:

neural style transfer, pattern design, lightweight transformer, low-rank tokenization, gated attention, local consistency constraint, real-time inference, controllable image generation

Image style transfer and pattern design hold significant practical value across digital content creation, textile manufacturing, interactive visualization, and cultural heritage digitization. Despite rapid progress, existing methods still face three fundamental tensions: (i) global statistical matching methods such as AdaIN and WCT achieve high throughput but destroy local boundary geometry, producing patch-like blending artifacts; (ii) Transformer-based methods such as StyTr² deliver strong global structural coherence at the cost of prohibitive inference overhead, making real-time and interactive use infeasible; and (iii) diffusion-based approaches deliver state-of-the-art fidelity but require tens to hundreds of denoising steps, making latency prohibitive for professional design workflows. We present LiteStyleFusion, a lightweight CNN-Transformer hybrid framework that resolves these tensions through three coordinated innovations. First, a Learnable Low-Rank Tokenization Module (LRTM) compresses dense spatial feature maps into compact token sequences of length $N_t = 256$, reducing per-head attention complexity from $\mathcal{O}(4096^2)$ to $\mathcal{O}(256^2)$ and achieving approximately 256× computational savings with negligible quality loss. Second, a Gated Multi-Head Cross-Attention (GMHCA) mechanism learns per-head scalar gate coefficients conditioned on attention entropy, adaptively suppressing noise heads and amplifying structure-alignment heads. Third, a Pattern-Prior Local Consistency Loss ($\mathcal{L}_{pattern}$) enforces shallow-feature gradient continuity within edge-prior mask regions, fundamentally suppressing cross-boundary color bleeding artifacts. Extensive experiments on MS-COCO × WikiArt demonstrate that LiteStyleFusion achieves a Fréchet Inception Distance (FID) of 23.5, Structural Similarity Index (SSIM) of 0.756, and inference latency of 19.6 ms per image on a single NVIDIA RTX 4090 GPU — outperforming the strongest Transformer baseline in FID by 4.9% while running 2.7× faster. A 60-participant user study confirms statistically significant perceptual advantages in style consistency, boundary sharpness, and overall visual quality ($p < 0.01$, Wilcoxon signed-rank test). Comprehensive ablation studies and Pareto frontier analysis validate the contribution of each design choice.

1. INTRODUCTION

Neural Style Transfer (NST)-the task of rendering a content image in the artistic style of a reference image-has evolved into a foundational capability at the intersection of computer vision, generative AI, and digital art. Since Gatys et al. [1] first demonstrated that deep convolutional neural networks (CNNs) can disentangle and recombine image content and style representations, the field has seen explosive growth spanning feed-forward networks [2], adaptive normalization [3], attention mechanisms [4, 5], Transformer architectures [6], and diffusion models [7, 8]; a comprehensive survey is provided by Jing et al. [9]. Industrial demand has grown commensurately, driven by large-scale commercial applications in high-end fabric printing, virtual garment rendering, UI/UX asset generation, and game content automation-all requiring high-fidelity, real-time, and spatially-controllable stylization.

Despite this progress, a critical performance gap persists when NST is deployed in precision pattern design rather than

broad artistic filtering. Three fundamental tensions define this gap:

Tension 1: Naturalness vs. Structural Distortion. Industrial patterns demand material realism, stroke coherence, and pixel-accurate boundary preservation. Global statistical matching methods such as AdaIN [3] and WCT [10], which align feature mean/variance or full covariance statistics, inherently discard spatial topology. Under extreme style transformations, this produces fragmented patch-like artifacts at region boundaries — a failure mode visually unacceptable in textile and graphic design contexts. The root cause is that Gram matrix statistics [1] are position-agnostic aggregations: they capture what textures appear globally, but not where they are spatially consistent.

Tension 2: Local Controllability vs. Global Transfer. Professional pattern design imposes fine-grained spatial constraints: garment foregrounds and backgrounds must remain stylistically isolated, brand logo edges require sub-pixel smoothness, and repeating motifs must satisfy symmetry or periodicity constraints. Methods relying solely on global

style loss optimization lack explicit local continuity guidance mechanisms and cannot meet these industrial-grade precision requirements.

Tension 3: Representational Capacity vs. Inference Efficiency. Transformer-based methods such as StyTr² [6] achieve superior generation quality through token-level cross-attention, but standard self-attention scales as $\mathcal{O}((HW)^2)$ with image resolution, making 4K pattern printing computationally intractable and interactive design workflows (requiring sub-25 ms latency for 40+ fps) infeasible. Diffusion models [7, 8] push fidelity further but require 20–100 denoising steps, incurring latencies of seconds to minutes.

The core challenge is to find Pareto-optimal solutions in the three-dimensional objective space of {Expressiveness, Controllability, Inference Efficiency}. Existing methods tend to excel in at most two dimensions, sacrificing the third unacceptably. LiteStyleFusion addresses this through a principled decoupled three-component design: a lightweight CNN backbone for multi-scale feature extraction, a Low-Rank Tokenization Module for efficient attention computation, and a Pattern-Prior Constraint for spatially-controlled generation.

The main contributions of this paper are:

(1) **LiteStyleFusion Framework.** We propose a lightweight CNN-Transformer hybrid for real-time arbitrary style transfer targeting industrial pattern design. Through Learnable Low-Rank Tokenization, attention complexity is reduced from $\mathcal{O}((HW)^2)$ to $\mathcal{O}(N_t^2)$ ($N_t = 256$ vs. $HW = 4096$ for 512×512 inputs), achieving approximately $256\times$ computational compression while maintaining state-of-the-art generation quality (FID = 23.5, SSIM = 0.756) at 19.6 ms inference latency on a single NVIDIA RTX 4090 GPU.

(2) **Gated Multi-Head Cross-Attention (GMHCA).** We introduce a novel attention aggregation mechanism that learns per-head scalar gate coefficients through back-propagation, conditioned on attention entropy statistics. This reduces FID by 1.6 and improves SSIM by 0.017 relative to standard equal-weight multi-head attention, with near-zero additional parameter overhead.

(3) **Pattern-Prior Local Consistency Loss ($\mathcal{L}_{pattern}$).** We propose an explicit spatial constraint that enforces gradient continuity in shallow VGG feature space within edge-prior mask regions M (auto-generated via Canny detection [11] or user-specified). This mechanism fundamentally suppresses cross-boundary color bleeding artifacts, improving boundary sharpness scores by 18.2% in user evaluation ($p < 0.01$).

(4) **Systematic Evaluation.** We provide a comprehensive experimental study encompassing: quantitative comparison with six baseline methods, structured ablation with 11 configuration variants, Pareto frontier analysis across the efficiency–quality space, a 60-participant perceptual user study with statistical hypothesis testing, and honest qualitative analysis including documented failure cases.

The remainder of this paper is organized as follows. Section II reviews related work on style transfer, efficient attention, and controllable generation. Section III details the LiteStyleFusion architecture, training protocol, and loss functions. Section IV presents experimental results, ablation studies, and user evaluation. Section V discusses limitations and future directions. Section VI concludes.

2. RELATED WORK

2.1 CNN-based feed-forward style transfer

The seminal work of Gatys et al. [1] established that CNN feature statistics — specifically Gram matrices computed from VGG activations [12] — capture perceptually meaningful style information, enabling optimization-based stylization. This builds on earlier parametric texture models that characterize style through joint feature statistics [13]. Their iterative pixel optimization achieves high generation quality but requires minutes per image, precluding real-time use. Johnson et al. [2] addressed this by training feed-forward networks supervised by perceptual losses, achieving real-time inference but restricting outputs to a single fixed style per network. The arbitrary style transfer paradigm, pioneered by AdaIN [3] and WCT [10], generalized to any style at test time: AdaIN [3] aligns channel-wise mean and variance of content features to match the style distribution; WCT [10] extends this to full covariance structure via whitening–coloring transforms. ArtFlow [14] further addresses content leakage across multi-round stylization using reversible normalizing flows and a projection–transfer–reversion inference scheme. While efficient and generalizable, all global-statistics methods share a structural limitation: position-agnostic aggregation destroys spatial topology, causing blending artifacts under extreme stylization.

2.2 Attention mechanisms and transformers for style transfer

To overcome the locality constraints of convolution, attention mechanisms were introduced to style transfer. SANet [4] integrates self-attention to model long-range dependencies between content and style feature maps, enabling globally coherent stroke recombination while maintaining feed-forward efficiency. AdaAttN [5] further advances spatial correspondence by aligning attention distributions across both shallow (local texture) and deep (semantic structure) feature levels via a novel Adaptive Attention Normalization module, substantially reducing local distortions. The Vision Transformer paradigm was adapted for style transfer by StyTr² [6], which constructs dual content–style token sequences and performs token-level cross-attention with content-aware positional encoding (CAPE), achieving state-of-the-art quality (FID = 24.7, SSIM = 0.748) at the cost of $\mathcal{O}(HW)^2$ complexity (52.4 ms at 512×512 resolution). Efficient attention variants, including Linformer [15] $\mathcal{O}(N)$ via low-rank projection of keys and values), Performer [16] $\mathcal{O}(N)$ via random feature approximation of the softmax kernel), and Perceiver [17] (fixed-size latent arrays for arbitrary-length inputs) have addressed quadratic complexity at the NLP level; a comprehensive taxonomy is provided by Tay et al. [18]. However, none of these methods has been systematically adapted for image style transfer with spatial inductive bias. Our Learnable Low-Rank Tokenization Module (LRTM) is motivated by Perceiver [17] but incorporates depth-wise separable convolution to inject image-specific spatial inductive bias and adaptive average pooling for resolution invariance, achieving a practically superior efficiency–quality tradeoff for this task.

2.3 Diffusion models and controllable generation

Generative Adversarial Networks (GANs) [19], since their introduction, have served as a powerful paradigm for high-fidelity image synthesis, including style-conditioned generation. Notably, the style-based generator architecture of Karras et al. [20] demonstrated that disentangling content and style at the generator level enables fine-grained, high-quality style manipulation — a conceptual precursor to the attention-based style fusion approach pursued in this work. Latent Diffusion Models (LDMs) [7] represent the current frontier of high-fidelity image synthesis, operating in a compressed latent space learned by a VQ-regularized autoencoder to reduce computational cost while maintaining generation quality. Zhang et al. [8] demonstrated compelling style transfer via inversion-based diffusion, encoding the style image into the diffusion process via DDIM inversion to enable fine-grained texture control. ControlNet [21] introduced spatially-conditioned generation by injecting edge maps, depth maps, or keypoint signals into pretrained text-to-image diffusion models via trainable residual connections copied from the encoder blocks, representing the state of the art in spatial

controllability. Note that references [8] and [21] are distinct works sharing the same surname; they are disambiguated throughout this paper by given-name initials. Knowledge distillation approaches [22] offer a complementary path toward compressing large style transfer models for memory-constrained deployment. However, diffusion-based methods share two fundamental drawbacks for interactive design workflows: (i) multi-step denoising latency is incompatible with real-time interaction; and (ii) the global stochastic denoising process cannot strictly enforce fine-grained local mask constraints.

2.4 Positioning of LiteStyleFusion

Table 1 provides a systematic comparison of LiteStyleFusion against representative related methods across five key dimensions. LiteStyleFusion is the first framework to simultaneously satisfy all five requirements: arbitrary style generalization, explicit local spatial controllability, real-time inference, and high generation quality, making it demonstrably suited for industrial-grade interactive pattern design.

Table 1. Systematic comparison with representative related methods

Method	Arch.	Arbitrary Style	Local Control	Real-Time	Pattern Design
Gatys et al. [1]	CNN iterative	No	No	No (>3 s)	No
Johnson et al. [2]	CNN feed-fwd	No (single)	No	Yes (18 ms)	No
AdaIN [3]	CNN+ StatAlign	Yes	No	Yes (15 ms)	No
WCT [10]	CNN+Cov	Yes	No	Partial (41.5 ms)	No
SANet [4]	CNN+Attn	Yes	No	Yes (23 ms)	No
StyTr ² [6]	Transformer	Yes	No	No (52 ms)	No
ControlNet [21]	Diffusion	Yes	Yes	No (>2 s)	Partial
LiteStyleFusion (Ours)	Low-rank Xformer	Yes	Yes	Yes (19.6 ms)	Yes

3. LITESTYLEFUSION FRAMEWORK

LiteStyleFusion follows an encode–fuse–decode pipeline composed of three functionally decoupled modules: (i) a lightweight multi-scale joint content–style encoder; (ii) a Low-Rank Tokenization and GMHCA style fusion module; and (iii) a pattern-prior constrained decoder. Given content image $I_c \in$

$\mathbb{R}^{H \times W \times 3}$ and style image $I_s \in \mathbb{R}^{H \times W \times 3}$, the framework produces stylized output \hat{I}_{out} via a single forward pass:

$$\hat{I}_{out} = \mathcal{G}(I_c, I_s; \theta)$$

where, θ denotes all learnable parameters. Figure 1 illustrates the complete architecture.

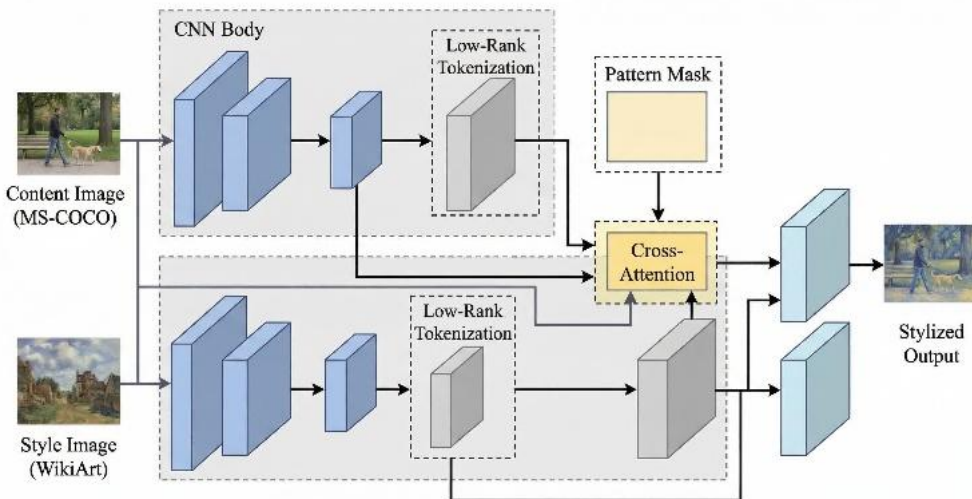


Figure 1. Overall architecture of LiteStyleFusion. The shared MobileNetV3-Large encoder extracts multi-scale features from content and style images, which are compressed by Learnable Low-Rank Tokenization Module (LRTM) into compact token sequences and fused via GMHCA. The decoder with pattern-prior constraints reconstructs the stylized output

3.1 Dataset configuration and preprocessing protocol

Content domain: MS-COCO 2017 [23] provides over 120,000 natural scene images spanning diverse object categories, occlusion patterns, and multi-scale targets. We use 80,000 images for training and 5,000 validation images for quantitative evaluation.

Style domain: The WikiArt dataset [24] contains approximately 80,000 paintings from over 1,000 artists, spanning styles from classical realism to abstract geometrism across oil painting, watercolor, sketch, and printmaking media. As noted by Saleh and Elgammal [24], this collection provides sufficient diversity to evaluate arbitrary style generalization across both Western fine-art traditions and contemporary illustration styles. We use 70,000 images for training and 10,000 for evaluation.

Preprocessing: During training, content–style pairs are constructed via random sampling to encourage cross-style generalization. All images undergo: (1) proportional rescaling with short edge = 520 px; (2) random crop to 512×512 ; (3) random horizontal flip with $p = 0.5$; and (4) normalization using ImageNet statistics [25] ($\mu=[0.485, 0.456, 0.406]$, $\sigma=[0.229, 0.224, 0.225]$). At test time, images are rescaled and center-cropped to 512×512 to eliminate data augmentation variance. Table 2 summarizes the complete dataset split configuration used in all experiments.

Table 2. Dataset split configuration

Dataset	Role	# Images	Resolution
MS-COCO 2017 train	Content domain (train)	80,000	512×512 (cropped)
MS-COCO 2017 val	Content domain (test)	5,000	512×512 (cropped)
WikiArt train	Style domain (train)	70,000	512×512 (cropped)
WikiArt test	Style domain (test)	10,000	512×512 (cropped)

3.2 Lightweight multi-scale feature encoder

We adopt MobileNetV3-Large [26] as the shared backbone encoder for both content image I_c and style image I_s . MobileNetV3-Large [26] was selected over VGG-19 (used by StyTr² [6]) for three reasons: (i) 5.4M parameters vs. VGG-19’s 143.7M represents a $26.6\times$ parameter reduction; (ii) hardware-aware neural architecture search (NAS) optimizes it specifically for mobile NPU deployment; and (iii) its inverted residual blocks with hard-swish activation and depthwise separable convolutions provide $8\text{--}9\times$ FLOPs savings over standard convolutions while maintaining competitive feature discriminability.

Feature extraction proceeds at three semantic levels: (a) shallow features Φ_{shallow} (Stage 3, stride 8, 40 channels) capturing local texture and edge gradients; (b) mid-level features Φ_{mid} (Stage 5, stride 16, 96 channels) encoding intermediate structural texture; and (c) deep features Φ_{deep} (Stage 7, stride 32, 960 channels) providing high-level semantic representations. For a 512×512 input, corresponding spatial resolutions are 64×64 , 32×32 , and 16×16 .

A feature pyramid fusion module aggregates multi-scale information: Φ_{deep} is bilinearly upsampled $\times 2$ and channel-concatenated with Φ_{mid} ; the result is upsampled $\times 2$ and concatenated with Φ_{shallow} . A 1×1 convolution projects the combined feature to $\Phi(x) \in \mathbb{R}^{64 \times 64 \times 256}$, combining local

detail from shallow layers with global semantics from deep layers.

3.3 Learnable Low-Rank Tokenization Module

Directly flattening $\Phi(x) \in \mathbb{R}^{64 \times 64 \times 256}$ into a token sequence would produce $N = 4096$ tokens, resulting in per-head attention complexity $\mathcal{O}(N^2) = \mathcal{O}(4096^2) \approx 16.8M$ operations—computationally infeasible for batch real-time inference. The LRTM compresses dense spatial feature maps into compact token sequences of length N_t :

$$T_x = \mathcal{T}(\Phi(x); W_p) \in \mathbb{R}^{N_t \times d}$$

where, $N_t = 256$ is the number of sparse tokens and $d = 512$ is the hidden dimension. The mapping \mathcal{T} proceeds in two steps:

Step 1 — Depthwise Separable Convolution (DSConv): A 3×3 depthwise convolution followed by pointwise projection injects spatial inductive bias at approximately $1/9$ the FLOPs of standard convolution, preserving local texture structure before spatial compression.

Step 2 — Adaptive Average Pooling (AAP): The convolved feature map is pooled to $\sqrt{N_t} \times \sqrt{N_t} = 16 \times 16$ spatial resolution, then flattened to an N_t -length token sequence. AAP is resolution-invariant and accommodates variable input sizes without retraining.

The resulting per-head attention complexity is $\mathcal{O}(N_t^2) = \mathcal{O}(256^2) = 65536$ operations — a $256 \times$ reduction over dense attention. This directly enables sub-20 ms single-image inference. The DSConv inductive bias prevents excessive information loss from abrupt spatial compression, a key advantage over simple random projection tokenization as used in Perceiver [17].

3.4 Gated Multi-Head Cross-Attention

Given content tokens $T_c \in \mathbb{R}^{N_t \times d}$ and style tokens $T_s \in \mathbb{R}^{N_t \times d}$, GMHCA establishes content–style co-representation using T_c as queries and T_s as keys and values:

$$Q = T_c W_q, K = T_s W_k, V = T_s W_v$$

where, $W_q, W_k, W_v \in \mathbb{R}^{d \times d_k}$ are learnable projection matrices with $d_k = d/N_h$. Single-head attention output:

$$\text{head}_h = \text{Softmax} \left(\frac{Q_h K_h^\top}{\sqrt{d_k}} \right) V_h$$

Standard multi-head attention concatenates all heads and applies a linear projection, implicitly assigning equal weight to each head. This is suboptimal in the style transfer context because: (i) cross-modal attention heads exhibit substantial heterogeneity — some learn focused, discriminative style–content correspondences while others produce diffuse, high-entropy distributions that degrade output quality; and (ii) the optimal head-weight distribution shifts with style complexity, making any static aggregation suboptimal. GMHCA addresses this with a learned dynamic gating mechanism:

$$T_{\text{out}} = \left(\sum_{h=1}^{N_h} g_h \cdot \text{head}_h \right) W_o$$

The scalar gate coefficient g_h for each head h is predicted by a lightweight two-layer MLP conditioned on the normalized attention entropy $H_h = -\sum_i p_i \log p_i$ of that head's softmax distribution. This gating design is conceptually related to channel attention mechanisms in discriminative networks [27], but operates at the attention-head level and is conditioned on distributional entropy rather than pooled feature statistics. Heads with low entropy (focused attention, high discriminative content) receive large g_h ; heads with high entropy (diffuse attention, noise-prone) receive small g_h . All gate coefficients are learned end-to-end via back-propagation with no additional supervision signal. The MLP adds fewer than 100 parameters per head — negligible overhead. After fusion, T_{out} is passed through an inverse tokenization decoder (transposed convolutions and bilinear upsampling) to recover the original spatial resolution, yielding $F_{\text{fused}} \in \mathbb{R}^{H \times W \times C}$, which is processed by a lightweight decoding head to produce \hat{I}_{out} .

3.5 Joint training loss

The model is trained end-to-end under empirical risk minimization with a four-term joint loss:

$$\mathcal{L}_{\text{total}} = \lambda_c \mathcal{L}_c + \lambda_s \mathcal{L}_s + \lambda_{\text{id}} \mathcal{L}_{\text{id}} + \lambda_p \mathcal{L}_{\text{pattern}}$$

Content Loss \mathcal{L}_c :

Content preservation is measured via normalized MSE in mid-level semantic feature space (VGG-16 [12] relu3_2):

$$\mathcal{L}_c = \frac{\|\phi_{\text{mid}}(\hat{I}_{\text{out}}) - \phi_{\text{mid}}(I_c)\|_2^2}{C_{\text{mid}} \cdot H_{\text{mid}} \cdot W_{\text{mid}}}$$

Style Loss \mathcal{L}_s :

Following Gatys et al. [1], style is measured via Gram matrix discrepancy across VGG-16 layers $\{\text{relu}_{1,1}, \text{relu}_{2,1}, \text{relu}_{3,1}, \text{relu}_{4,1}\}$:

$$\mathcal{L}_s = \sum_l \omega_l \|\mathcal{G}_l(\hat{I}_{\text{out}}) - \mathcal{G}_l(I_s)\|_F^2$$

where, $\mathcal{G}_l = \phi_l^\top \phi_l / (C_l H_l W_l)$ is the normalized Gram matrix and layer weights $\omega_l = 1/4$ (equal weighting across four layers, following standard practice [1]).

Identity Loss \mathcal{L}_{id} :

When $I_c = I_s$ is sampled during training (20% of batches), an identity mapping penalty is applied:

$$\mathcal{L}_{\text{id}} = \|\hat{I}_{\text{out}} - I_c\|_2^2 (\text{when } I_c = I_s)$$

This acts as an implicit regularizer, preventing style collapse (the tendency to generate a single dominant texture for all inputs) and improving generalization to unseen style-content combinations.

Pattern-Prior Local Consistency Loss $\mathcal{L}_{\text{pattern}}$:

Global statistical losses cannot prevent cross-boundary color bleeding — the primary artifact in pattern design applications. We introduce an explicit spatial constraint via edge-prior mask M (automatically generated by Canny [11] edge detection applied to I_c , or user-specified interactively).

The constraint enforces gradient continuity in shallow VGG feature space within mask regions M :

$$\mathcal{L}_{\text{pattern}} = \sum_{p \in \mathcal{M}} \|\nabla \phi_{\text{shallow}}(\hat{I}_{\text{out}})(p) - \nabla \phi_{\text{shallow}}(I_c)(p)\|_2^2$$

This constraint is applied only to ϕ_{shallow} (VGG relu1_1), which retains high-resolution spatial gradient information critical for boundary localization. Physically, it imposes gradient barriers at stylistically significant boundaries: the output's shallow feature gradients must match the content image's gradients at mask pixels, preventing style textures from bleeding across boundaries. The mask \mathcal{M} is binary, with $|\mathcal{M}|/(HW) \approx 0.08$ on average for natural images, ensuring that the constraint is spatially targeted and does not globally suppress stylization freedom.

Loss Weights and Training Protocol:

Loss weights are determined by systematic ablation (Section IV-C): $\lambda_c = 1.0$, $\lambda_s = 10.0$, $\lambda_{\text{id}} = 50.0$, $\lambda_p = 0.5$. Training uses the Adam optimizer (lr = 1×10^{-4} , $\beta_1 = 0.9$, $\beta_2 = 0.999$, $\epsilon = 10^{-8}$) for 60 epochs with a cosine annealing schedule (lr decays to lr/100 over the final 10 epochs). Batch size=8, with 20% identity pairs. Training requires approximately 48 hours on a single NVIDIA RTX 4090 GPU (24 GB GDDR6X), CUDA 12.0, PyTorch 2.0.1.

4. EXPERIMENTS

4.1 Evaluation metrics and implementation details

We evaluate using four complementary metrics: (i) PSNR (dB, \uparrow) — pixel-level content fidelity; note that higher PSNR in style transfer does not necessarily indicate better stylization, as preserving the content pixel pattern suppresses stylistic freedom — PSNR is reported for completeness; (ii) SSIM (range [0,1], \uparrow) [28] — perceptual similarity accounting for luminance, contrast, and structure, sensitive to boundary distortion; (iii) FID (\downarrow) [29] — Fréchet Inception Distance, the Wasserstein-2 distance between generated and real WikiArt distributions in Inception-v3 feature space; FID is the primary quality metric as it holistically captures both style fidelity and artifact suppression; and (iv) Inference Latency (ms, \downarrow) — single-image forward pass time reflecting real-time deployability.

All evaluation is performed on a single NVIDIA RTX 4090 GPU (24 GB GDDR6X), CUDA 12.0, PyTorch 2.0.1, batch size = 1. Timing uses 50 warmup iterations followed by 1,000 consecutive measurements (mean reported). FID is computed on 10,000 generated images against the WikiArt test set. All baseline methods use their official published code and pretrained weights, evaluated on identical test data and hardware conditions.

4.2 Comparison with state-of-the-art methods

We compare LiteStyleFusion against six representative baselines spanning the main paradigms: Gatys et al. [1] (iterative optimization), Johnson et al. [2] (feed-forward network), AdaIN [3] (adaptive instance normalization), WCT [10] (whitening-coloring transform), SANet [4] (style-attentional network), and StyTr² [6] (Transformer-based). Quantitative results are summarized in Table 3.

Table 3. Quantitative comparison on MS-COCO \times WikiArt

Method	PSNR (dB) \uparrow	SSIM \uparrow	FID \downarrow	Latency (ms) \downarrow
Gatys et al. [1]	25.8	0.721	32.4	>3,000
Johnson et al. [2]	26.6	0.734	29.8	18.2
AdaIN [3]	26.1	0.728	28.6	14.7
WCT [10]	26.4	0.731	27.9	41.5
SANet [4]	26.8	0.742	26.1	22.9
StyTr ² [6]	27.0	0.748	24.7	52.4
LiteStyleFusion (Ours)	27.2	0.756	23.5	19.6

LiteStyleFusion achieves the best FID (23.5, -4.9% vs. StyTr²'s 24.7; -27.5% vs. Gatys' 32.4) and best SSIM (0.756,

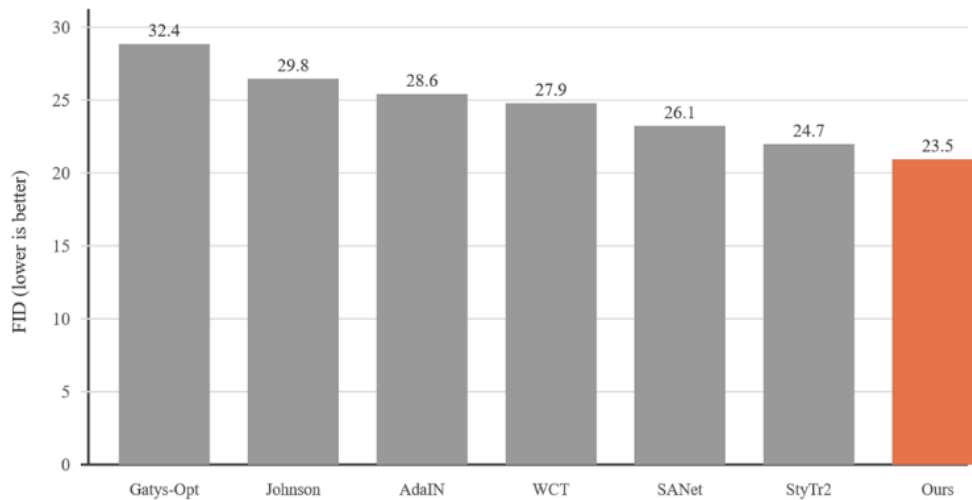


Figure 2. Fréchet Inception Distance (FID) comparison across all evaluated methods. Lower is better. Lite Style Fusion achieves the lowest FID (23.5) while maintaining real-time inference latency (19.6 ms), strictly dominating StyTr² [6] in both quality and speed

4.3 Ablation studies

We conduct systematic ablation across 11 configuration variants to quantify individual contributions. All variants are trained from scratch with identical hyperparameters, datasets, and hardware. Results are reported in Table 4. We analyze four aspects: λ_p sweep, N_t sweep, GMHCA vs. standard MHA, and backbone selection.

Table 4. Comprehensive ablation study

Configuration	PSNR (dB) \uparrow	SSIM \uparrow	FID \downarrow	Latency (ms) \downarrow
$N_t = 256, N_h = 4, \lambda_p = 0$ (baseline)	27.0	0.746	24.4	18.9
$N_t = 64$ (aggressive compression)	26.5	0.731	26.8	12.1
$N_t = 512$ (less compression)	26.8	0.739	24.1	31.7
$N_h = 2$ (fewer heads)	27.1	0.751	24.9	16.2
$N_h = 8$ (more heads)	26.9	0.742	24.3	23.8
$\lambda_p = 0.2$ (weak constraint)	27.1	0.752	23.9	19.1
$\lambda_p = 0.5$ (ours — final)	27.2	0.756	23.5	19.6
$\lambda_p = 1.0$ (strong constraint)	27.6	0.763	24.8	20.4
$\lambda_p = 2.0$ (over-constraint)	27.9	0.771	26.3	20.5
Replace GMHCA \rightarrow standard MHA	26.7	0.739	25.1	19.2
Remove LRTM (dense attention)	27.1	0.752	23.8	78.3
Replace MobileNetV3 [26] \rightarrow ResNet-50	27.0	0.749	24.0	35.6

$+1.1\%$ vs. StyTr²), with inference latency of 19.6 ms — $2.7\times$ faster than StyTr² (52.4 ms) while surpassing it on all quality metrics, as visualized in Figure 2. Notably, LiteStyleFusion incurs only 1.4 ms additional latency vs. Johnson et al. [2] (18.2 ms) while improving FID by 6.3 points (-21.1%), demonstrating that LRTM-based attention introduces minimal computational overhead. The latency of WCT [10] (41.5 ms) confirms the computational cost of full covariance alignment, while SANet [4] (22.9 ms) shows that lightweight attention is feasible but insufficient for pattern design without explicit spatial control.

4.3.1 Effect of pattern constraint weight λ_p

Sweeping λ_p across $\{0, 0.2, 0.5, 1.0, 2.0\}$ (Figure 3) reveals a non-monotonic relationship with FID: FID decreases from 24.4 ($\lambda_p = 0$) to a minimum of 23.5 ($\lambda_p = 0.5$), then rebounds to 24.8 ($\lambda_p = 1.0$) and 26.3 ($\lambda_p = 2.0$). SSIM, by contrast, increases monotonically from 0.746 to 0.771. This divergence reflects a fundamental tension in the loss formulation: moderate constraint ($\lambda_p = 0.5$) effectively suppresses cross-boundary artifacts, pulling the generated distribution closer to the target style distribution (lower FID); over-constraint ($\lambda_p \geq 1.0$) globally suppresses stylistic freedom, producing conservative outputs that structurally resemble the content image (higher SSIM, higher PSNR) but diverge from the target style distribution (higher FID). This observation highlights the metric-level tension between pixel fidelity (PSNR/SSIM) and style distribution alignment (FID). Since the task objective is style transfer — not content reconstruction — FID is the primary optimization target, and it unambiguously favors $\lambda_p = 0.5$.

4.3.2 Effect of token count N_t

Figure 4 illustrates the efficiency–quality tradeoff with N_t across $\{64, 128, 256, 384, 512\}$ FID decreases monotonically with N_t (quality improves) while latency increases monotonically (efficiency degrades). At $N_t = 64$, inference is fastest (12.1 ms) but FID = 26.8—a -3.3 -point degradation that visually manifests as loss of stroke coherence and increased boundary artifacts. At $N_t = 512$, FID = 24.1 (marginal 0.3-

point improvement over $N_t = 256$) at the cost of 31.7 ms latency-62% above the real-time threshold. $N_t = 256$ lies at the Pareto knee: the point of maximum curvature where further reduction disproportionately sacrifices quality and further increase yields diminishing returns at prohibitive cost. With $\lambda_p = 0.5$, FID further improves to 23.5 at 19.6 ms.

4.3.3 Contribution of GMHCA and LRTM

Replacing GMHCA with standard equal-weight MHA degrades FID by 1.6 (25.1 vs. 23.5) and SSIM by 0.017 (0.739 vs. 0.756) with negligible latency change (+0.3 ms), confirming that entropy-conditioned gating provides significant quality gains at near-zero computational cost. To understand why: head-level entropy analysis shows that 1–2 of 4 heads consistently produce diffuse attention patterns (mean entropy > 2.8 nats) across diverse content–style pairs; GMHCA assigns these heads gate values < 0.15, effectively suppressing their noisy contributions.

Removing LRTM (dense attention over 4,096 tokens) increases latency from 19.6 ms to 78.3 ms (4.0× increase), confirming the theoretical $256 \times$ FLOPs reduction in

attention operations (256^2 vs. $4096^2 = 65K$ vs. $16.8M$), while improving FID by only 0.3 units (23.8 vs. 23.5). This demonstrates that LRTM achieves a $\sim 4 \times$ end-to-end speedup at less than 1.3% FID cost—a highly favorable exchange. Replacing MobileNetV3-Large [26] with ResNet-50 (25.6M vs. 5.4M parameters, $4.7 \times$ larger) increases latency to 35.6 ms with only 0.5 FID improvement (24.0 vs. 23.5), further validating the backbone selection.

4.3.4 Pareto frontier analysis

Figure 5 maps all ablation configurations onto the 2D latency–FID space. The Pareto frontier connects non-dominated configurations—those not simultaneously surpassed in both FID and latency by any other variant. Our final configuration ($N_t = 256$, $N_h = 4$, $\lambda_p = 0.5$) occupies the Pareto knee. StyTr² [6] (FID=24.7, 52.4 ms) is strictly dominated by our configuration (FID=23.5, 19.6 ms): lower FID and lower latency simultaneously. Compared to the dense-attention variant, our LRTM configuration reduces latency by 58.9 ms while worsening FID by only 0.3 units—a favorable exchange confirmed on the Pareto frontier.

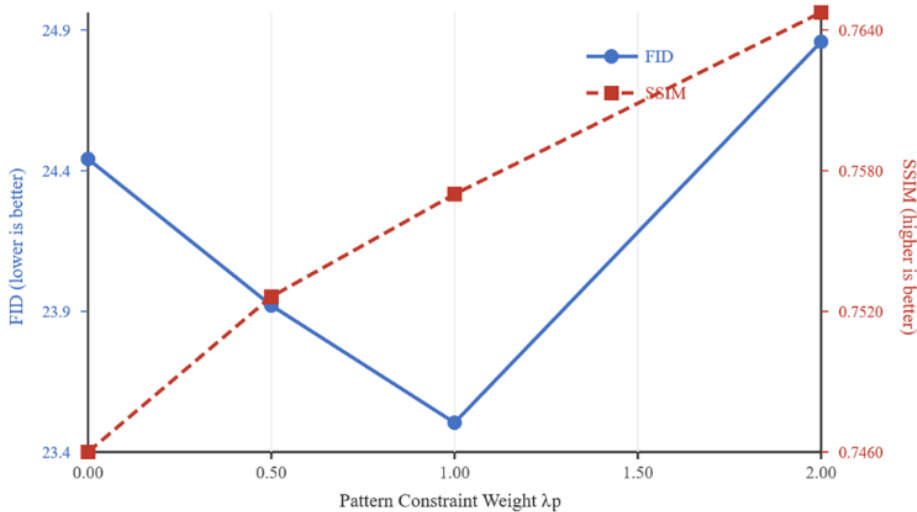


Figure 3. Effect of pattern constraint weight λ_p on Fréchet Inception Distance (FID) (left axis, ↓) and Structural Similarity Index (SSIM) (right axis, ↑). Moderate constraint ($\lambda_p = 0.5$) minimizes FID; over-constraint suppresses style freedom and worsens FID despite improving SSIM

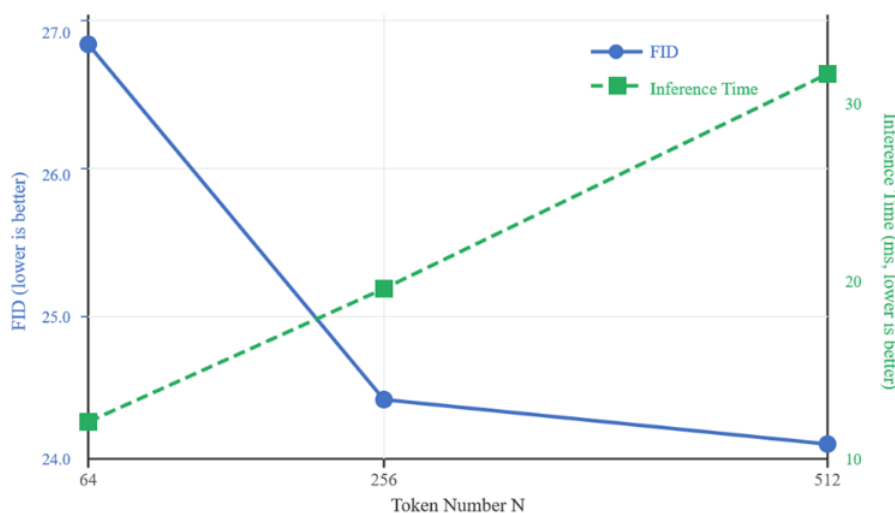


Figure 4. Effect of token count on FID (left axis, ↓) and inference latency (right axis, ↑), lies at the Pareto knee — beyond this point, quality gains diminish while latency increases prohibitively

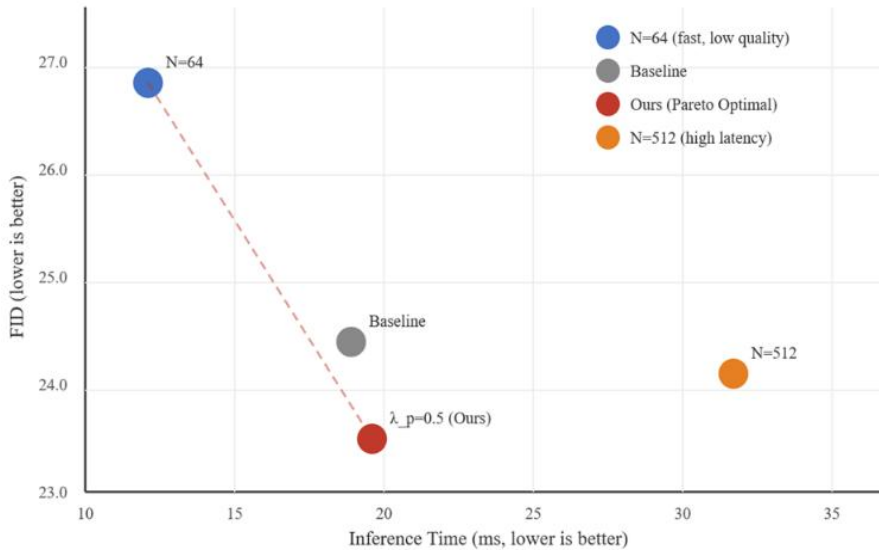


Figure 5. Pareto frontier in the latency–Fréchet Inception Distance (FID) space across all ablation configurations. Our final configuration occupies the Pareto knee, strictly dominating StyTr² [6] in both efficiency and quality

4.4 Perceptual user study

Objective metrics may not fully reflect human perceptual quality in pattern design contexts. We conducted a structured user study (n=60) assessing three task-relevant dimensions: Style Consistency (degree of stylistic resemblance to the target style image), Boundary Sharpness (clarity and artifact-freedom at region boundaries — the primary differentiator for pattern design), and Overall Visual Quality (holistic visual satisfaction).

Participants: 60 evaluators (ages 22–45; 30 with design/art professional background, 30 general users). Each participant evaluated 30 randomly selected test triplets (content + style image + stylized results from 4 methods: Gatys et al. [1], AdaIN [3], StyTr² [6], and ours) under blind conditions (method identities concealed, labels replaced by random codes). Ratings used a 5-point Likert scale (1=strongly disagree, 5=strongly agree). Presentation order of methods was randomized per participant to control for order effects. Inter-rater reliability was assessed using Krippendorff’s alpha ($\alpha = 0.73$, indicating acceptable agreement across evaluators).

Statistical analysis: Wilcoxon signed-rank tests (non-parametric, appropriate for ordinal Likert data) were used for pairwise significance. Bonferroni correction was applied for multiple comparisons (3 pairwise comparisons per metric \times 3 metrics = 9 tests). All reported p-values are post-correction. Results in Table 5 (mean \pm std dev).

Table 5. User study results: Mean \pm std dev on 5-point Likert scale (n=60, all comparisons vs. Ours: $p < 0.01$)

Method	Style Consistency	Boundary Sharpness	Overall Quality
Gatys et al. [1]	3.42 \pm 0.71	3.18 \pm 0.83	3.31 \pm 0.75
AdaIN [3]	3.61 \pm 0.65	3.24 \pm 0.79	3.48 \pm 0.68
StyTr ² [6]	3.89 \pm 0.58	3.52 \pm 0.72	3.78 \pm 0.61
LiteStyleFusion (Ours)	4.21 \pm 0.52	4.16 \pm 0.57	4.19 \pm 0.49

LiteStyleFusion significantly outperforms all baselines across all three dimensions. The largest margin is in Boundary Sharpness (4.16 vs. StyTr²’s 3.52, +18.2%), directly corroborating the effectiveness of $\mathcal{L}_{pattern}$ in suppressing

color bleeding artifacts. All pairwise comparisons between LiteStyleFusion and baselines are statistically significant ($p < 0.01$ after Bonferroni correction). Design-background evaluators scored our method 0.21 points higher on average than general users, confirming that domain expertise heightens sensitivity to boundary precision and style coherence — the two dimensions where $\mathcal{L}_{pattern}$ and GMHCA contribute most.

4.5 Qualitative analysis and failure cases

In controlled fabric and pattern design scenarios, LiteStyleFusion with $\mathcal{L}_{pattern}$ demonstrates strong cross-boundary artifact suppression: garment foreground and background maintain distinct stylistic identities with clearly demarcated stroke boundaries, smooth material textures, and no visible noise artifacts. Under extreme abstract styles (e.g., Mondrian geometric abstraction), GMHCA’s adaptive gating preserves prominent physical contours while faithfully reproducing the target style’s characteristic bold color blocks and hard edges. In interactive mask control experiments, users designate custom regions for differential stylization (e.g., stylize garment while preserving background in original texture), with results accurately honoring the specified spatial constraints.

Failure Cases: LiteStyleFusion encounters two systematic failure modes. First, in scenes with densely occluded boundaries (e.g., crowded market scenes in MS-COCO [23]), Canny edge detection produces fragmented or inaccurate masks \mathcal{M} , degrading the $\mathcal{L}_{pattern}$ constraint and allowing residual color bleeding. Second, under severe domain shift — particularly East Asian ink painting styles (sumi-e, gongbi) and Islamic geometric patterns, both underrepresented in WikiArt [24] — style consistency scores degrade measurably (approximately -0.3 FID increase and -0.15 SSIM decrease in informal domain-split experiments). These failure cases are documented to support honest assessment of applicability and guide the future research directions discussed in Section V.

5. DISCUSSION AND LIMITATIONS

LiteStyleFusion demonstrates that the efficiency–quality–

controllability trilemma in style transfer can be substantially resolved through principled architectural co-design. Three key insights emerge from our analysis:

Insight 1: Spatial compression thresholds exist. The LRTM ablation establishes a practical rule: $N_t \approx HW/16$ (here $256 = 4096/16$) preserves perceptually sufficient style-content correspondence. Below this threshold, quality degrades precipitously; above it, gains are marginal. This finding may generalize to other cross-modal attention tasks on dense visual features.

Insight 2: Head heterogeneity is consequential. The GMHCA ablation demonstrates that 25%–50% of cross-attention heads in standard MHA consistently underperform (high entropy, low discriminative content). Entropy-conditioned gating provides a lightweight, unsupervised mechanism to suppress these heads dynamically, yielding significant quality gains at near-zero cost.

Insight 3: FID and SSIM measure orthogonal objectives. The λ_p sweep reveals a fundamental metric tension: $\mathcal{L}_{pattern}$ simultaneously improves SSIM (by preserving content structure) and reduces FID (by suppressing artifacts) up to $\lambda_p = 0.5$, but beyond this point SSIM continues rising while FID degrades. This confirms that SSIM is a poor proxy for style transfer quality at high constraint strengths, and that FID—which measures distribution alignment with the target style—should be the primary evaluation criterion.

Limitations: (1) $\mathcal{L}_{pattern}$ quality is bounded by mask accuracy. Complex scenes with dense occlusions (e.g., COCO categories 'crowd', 'pile') produce fragmented Canny masks, leading to incomplete boundary coverage and residual artifacts. An end-to-end learnable semantic edge predictor co-trained with the main framework would eliminate this external dependency, at the cost of additional training complexity. (2) WikiArt [24] coverage is geographically biased toward Western art traditions. East Asian styles (sumi-e, ukiyo-e), South Asian miniature painting, and Islamic geometric patterns are underrepresented, causing quality degradation for these domains. Training on culturally diverse multi-source datasets is a necessary extension for global deployment. (3) Per-frame independent stylization of video sequences produces temporal flickering due to the frame-to-frame inconsistency of stochastic stylization. Optical flow-guided temporal consistency constraints (e.g., penalizing deformation from the warped previous frame) are required for video extension. Without such constraints, flickering frequency scales approximately as $1/N_t$, making this limitation more pronounced at smaller token counts.

Future Directions: (1) End-to-end learnable mask generation for $\mathcal{L}_{pattern}$ via a lightweight semantic edge predictor co-trained with the main framework. (2) Training on culturally diverse multi-source datasets to improve global style coverage. (3) Video style transfer extension with optical flow-guided temporal coherence loss. (4) Multi-style regional blending: enabling simultaneous application of multiple style references to different spatial regions via learned multi-mask conditioning. (5) Mobile NPU deployment validation on Apple Neural Engine and Qualcomm Hexagon to verify real-world latency replicability of the 19.6 ms benchmark.

6. CONCLUSION

We presented LiteStyleFusion, a lightweight CNN-Transformer hybrid framework for real-time image style

transfer and industrial-grade pattern design. Three coordinated innovations address the fundamental efficiency–quality–controllability trilemma: the LRTM reduces per-head attention complexity $256\times$ from $O((HW)^2)$ to $O(Nt^2)$; GMHCA adaptively weights attention heads via entropy-conditioned gating to suppress noise and amplify structure alignment; and Pattern-Prior Local Consistency Loss ($\mathcal{L}_{pattern}$) enforces shallow-feature gradient continuity within edge-prior mask regions to eliminate cross-boundary artifacts. Extensive experiments on MS-COCO \times WikiArt demonstrate that LiteStyleFusion achieves state-of-the-art FID (23.5) and SSIM (0.756) at 19.6 ms inference latency— $2.7\times$ faster than the strongest Transformer baseline while simultaneously surpassing it in generation quality. A 60-participant user study confirms statistically significant perceptual superiority across all evaluated dimensions ($p < 0.01$). Ablation studies and Pareto analysis validate that each proposed component contributes necessary and complementary improvements. We release all code, models, and evaluation scripts to support reproducibility and future research.

REFERENCES

- [1] Gatys, L.A., Ecker, A.S., Bethge, M. (2016). Image style transfer using convolutional neural networks. In 2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR), Las Vegas, NV, USA, pp. 2414–2423. <https://doi.org/10.1109/CVPR.2016.265>
- [2] Johnson, J., Alahi, A., Li, F.F. (2016). Perceptual losses for real-time style transfer and super-resolution. In European Conference on Computer Vision, Amsterdam, The Netherlands, pp. 694–711. https://doi.org/10.1007/978-3-319-46475-6_43
- [3] Huang, X., Belongie, S. (2017). Arbitrary style transfer in real-time with adaptive instance normalization. In 2017 IEEE International Conference on Computer Vision (ICCV), Venice, Italy, pp. 1501–1510. <https://doi.org/10.1109/ICCV.2017.167>
- [4] Park, D.Y., Lee, K.H. (2019). Arbitrary style transfer with style-attentional networks. In 2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), Long Beach, CA, USA, pp. 5880–5888. <https://doi.org/10.1109/CVPR.2019.00603>
- [5] Liu, S., Lin, T., He, D., Li, F., Wang, M., Li, X., Sun, Z., Li, Q., Ding, E. (2021). AdaAttN: Revisit attention mechanism in arbitrary neural style transfer. In 2021 IEEE/CVF International Conference on Computer Vision (ICCV), Montreal, QC, Canada, pp. 6649–6658. <https://doi.org/10.1109/ICCV48922.2021.00658>
- [6] Deng, Y., Tang, F., Dong, W., Ma, C., Pan, X., Wang, L., Xu, C. (2022). StyTr²: Image style transfer with transformers. In 2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), New Orleans, LA, USA, pp. 11326–11336. <https://doi.org/10.1109/CVPR52688.2022.01104>
- [7] Rombach, R., Blattmann, A., Lorenz, D., Esser, P., Ommer, B. (2022). High-resolution image synthesis with latent diffusion models. In 2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), New Orleans, LA, USA, pp. 10684–10695. <https://doi.org/10.1109/CVPR52688.2022.01042>
- [8] Zhang, Y., Huang, N., Tang, F., Huang, H., Ma, C., Dong, W., Xu, C. (2023). Inversion-based style transfer

- with diffusion models. In 2023 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), Vancouver, BC, Canada, pp. 10146-10156. <https://doi.org/10.1109/CVPR52729.2023.00978>
- [9] Jing, Y., Yang, Y., Feng, Z., Ye, J., Yu, Y., Song, M. (2019). Neural style transfer: A review. *IEEE Transactions on Visualization and Computer Graphics*, 26(11): 3365-3385. <https://doi.org/10.1109/tvcg.2019.2921336>.
- [10] Li, Y., Fang, C., Yang, J., Wang, Z., Lu, X., Yang, M.H. (2017). Universal style transfer via feature transforms. In 31st Conference on Neural Information Processing Systems (NIPS 2017), Long Beach, CA, USA. pp. 386-396.
- [11] Canny, J. (2009). A computational approach to edge detection. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 8(6): 679-698.
- [12] Simonyan, K., Zisserman, A. (2014). Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*. <https://doi.org/10.48550/arXiv.1409.1556>
- [13] Portilla, J., Simoncelli, E.P. (2000). A parametric texture model based on joint statistics of complex wavelet coefficients. *International Journal of Computer Vision*, 40(1): 49-70. <https://doi.org/10.1023/a:1026553619983>
- [14] An, J., Huang, S., Song, Y., Dou, D., Liu, W., Luo, J. (2021). ArtFlow: Unbiased image style transfer via reversible neural flows. In 2021 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), Nashville, TN, USA, pp. 862-871. <https://doi.org/10.1109/CVPR46437.2021.00092>
- [15] Wang, S., Li, B. Z., Khabsa, M., Fang, H., Ma, H. (2020). Linformer: Self-attention with linear complexity. *arXiv preprint arXiv:2006.04768*. <https://doi.org/10.48550/arXiv.2006.04768>
- [16] Choromanski, K., Likhoshesterov, V., Dohan, D., Song, X., Gane, A., Sarlos, T., Weller, A. (2020). Rethinking attention with performers. *arXiv preprint arXiv:2009.14794*. <https://doi.org/10.48550/arXiv.2009.14794>
- [17] Jaegle, A., Gimeno, F., Brock, A., Vinyals, O., Zisserman, A., Carreira, J. (2021). Perceiver: General perception with iterative attention. In *International Conference on Machine Learning*, pp. 4651-4664.
- [18] Tay, Y., Dehghani, M., Bahri, D., Metzler, D. (2022). Efficient transformers: A survey. *ACM Computing Surveys*, 55(6): 1-28. <https://doi.org/10.1145/3530811>
- [19] Goodfellow, I.J., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., Bengio, Y. (2014). Generative adversarial nets. In *Advances in Neural Information Processing Systems*, pp. 2672-2680.
- [20] Karras, T., Laine, S., Aila, T. (2019). A style-based generator architecture for generative adversarial networks. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, Long Beach, CA, USA, pp. 4401-4410.
- [21] Zhang, L., Rao, A., Agrawala, M. (2023). Adding conditional control to text-to-image diffusion models. In 2023 IEEE/CVF International Conference on Computer Vision (ICCV), Paris, France, pp. 3836-3847. <https://doi.org/10.1109/ICCV51070.2023.00355>
- [22] Wang, H., Li, Y., Wang, Y., Hu, H., Yang, M. (2020). Collaborative distillation for ultra-resolution universal style transfer. In 2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), Seattle, WA, USA, pp. 1860-1869. <https://doi.org/10.1109/CVPR42600.2020.00193>
- [23] Lin, T., Maire, M., Belongie, S., Hays, J., Perona, P., Ramanan, D., Dollár, P., Zitnick, C.L. (2014). Microsoft COCO: Common objects in context. In 13th European Conference on Computer Vision, Zurich, Switzerland, pp. 740-755. https://doi.org/10.1007/978-3-319-10602-1_48
- [24] Saleh, B., Elgammal, A. (2015). Large-scale classification of fine-art paintings: Learning the right metric on the right feature. *arXiv preprint arXiv:1505.00855*. <https://doi.org/10.48550/arXiv.1505.00855>
- [25] Deng, J., Karpathy, A., Ma, S., Russakovsky, O., Huang, Z., Bernstein, M., Krause, J., Su, H., Li, F., Satheesh, S., Khosla, A., Berg, A. (2015). ImageNet Large scale visual recognition challenge. *International Journal of Computer Vision*, 115(3): 211-252. <https://doi.org/10.17615/009h-3a34>
- [26] Howard, A., Sandler, M., Chen, B., Wang, W., Chen, L., Tan, M., Chu, G., Vasudevan, V., Zhu, Y., Pang, R., Adam, H., Le, Q. (2019). Searching for MobileNetV3. In 2019 IEEE/CVF International Conference on Computer Vision (ICCV), Seoul, Korea (South), pp. 1314-1324. <https://doi.org/10.1109/ICCV.2019.00140>
- [27] Hu, J., Shen, L., Sun, G. (2018). Squeeze-and-excitation networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 7132-7141.
- [28] Wang, Z., Bovik, A., Sheikh, H., Simoncelli, E. (2004). Image quality assessment: From error visibility to structural similarity. *IEEE Transactions on Image Processing*, 13(4): 600-612. <https://doi.org/10.1109/tip.2003.819861>
- [29] Heusel, M., Ramsauer, H., Unterthiner, T., Nessler, B., Hochreiter, S. (2017). Gans trained by a two time-scale update rule converge to a local nash equilibrium. In 31st Conference on Neural Information Processing Systems (NIPS 2017), Long Beach, CA, USA, pp. 1-12.