

[Download book PDF](#)[Download book EPUB !\[\]\(666e09182d4cd268646ea700ea60dcdf\_img.jpg\)](#)

European Conference on Computer Vision

ECCV 2022: **Computer Vision – ECCV 2022** pp 717–734

## Language-Driven Artistic Style Transfer

[Tsu-Jui Fu](#) , [Xin Eric Wang](#) & [William Yang Wang](#)

Conference paper | [First Online: 29 October 2022](#)

**242** Accesses

Part of the [Lecture Notes in Computer Science](#) book series (LNCS, volume 13696)

### Abstract

---

Despite having promising results, style transfer, which requires preparing style images in advance, may result in lack of creativity and accessibility. Following human instruction, on the other hand, is the most natural way to perform artistic style transfer that can significantly improve controllability for visual effect applications. We introduce a new task—language-driven artistic style transfer (LDAST)—to manipulate the style of a content image, guided by a text. We propose contrastive language visual artist

(CLVA) that learns to extract visual semantics from style instructions and accomplish LDAST by the patch-wise style discriminator. The discriminator considers the correlation between language and patches of style images or transferred results to jointly embed style instructions. CLVA further compares contrastive pairs of content images and style instructions to improve the mutual relativeness. The results from the same content image can preserve consistent content structures. Besides, they should present analogous style patterns from style instructions that contain similar visual semantics. The experiments show that our CLVA is effective and achieves superb transferred results on LDAST.

---

Access provided by University of California, Santa  
Barbara

[Download](#) conference paper PDF

---

## 1 Introduction

---

Style transfer [14, 20, 21, 27, 28, 35] adopts appearances and visual patterns from another reference style images to manipulate a content image. Artistic style transfer has a considerable application value for creative visual design, such as image stylization and video effect [13, 19, 45, 59]. However, it requires preparing collections of style image in advance. It even needs to redraw new

references first if there is no expected style images, which is impractical due to an additional overhead. In contrast, language is the most natural way for humans to communicate. If a system can follow textual descriptions and automatically perform style transfer, we can significantly improve accessibility and controllability.

In this paper, we introduce Language-driven Artistic Style Transfer (LDAST). As illustrated in Fig. 1, LDAST treats a content image and a text as the input, and the style transferred result is manipulated based on the style description. It should preserve the structure of the content yet simultaneously modifies the style pattern that corresponds to the instruction. LDAST is different from the general language-based image-editing (LBIE) [9, 26, 31, 33] that aims at altering objects or properties of objects. The main challenge of LDAST is to extract visual semantics from language. Humans use not only explicit visual attributes but also visual content or emotional effects to describe style feelings. For example, it requires connecting *“water, sketching, and painting”* or *“peaceful, feel content”* with their visual concepts and further carrying out correlated style transfer.

We present contrastive language visual artist (CLVA), including language visual artist (LVA) and contrastive reasoning (CR), to perform style transfer conditioning on guided texts. LVA preserves content

structures from content images  $\mathcal{C}$  and extracts visual semantics from style instructions  $\mathcal{X}$ . LVA learns the latent style pattern based on the distinguishment between patches of style images or transferred results from the patch-wise style discriminator. Furthermore, CR boosts by comparing contrastive pairs where relative content images or style instructions should present similar content structures or style patterns.

To evaluate LDAST, we conduct experiments upon DTD<sup>2</sup> [50] and ArtEmis [1]. DTD<sup>2</sup> provides texture images with its colors or texture patterns in text. ArtEmis collects explanations of visual contents and emotional effects for artworks. We treat these annotations as style instructions for the challenging LDAST concerning visual attributes or human style feelings. The experiments show that our CLVA is effective for LDAST and achieves superb yet efficient transferred results on both automatic metrics and human evaluation. Our contributions are four-fold:

- We introduce LDAST that follows natural language for artistic style transfer;
- We present CLVA, which learns to extract explicit visual semantics from style instructions and provide sufficient style patterns for LDAST;
- We conduct the evaluation on DTD<sup>2</sup> and ArtEmis to consider diverse style instructions with visual attributes and emotional effects;



- Extensive experiments and qualitative examples demonstrate that our CLVA outperforms baselines regarding both effectiveness and efficiency.

**Fig. 1.**



Language-driven Artistic Style Transfer (LDAST). LDAST performs style transfer for a content image  $C$ , guided by the visual attribute (the lower row) or even the visual content and emotional effect (the upper row) from a style instruction  $\mathcal{X}$ .

## 2 Related Work

**Artistic Style Transfer.** Style transfer [6, 14, 16, 21, 22, 43, 47] redraws an image with a specific style. Since being a popular form of art, incorporating painting with digital design can produce attractive visual effect (VFX). In general, style transfer can be divided into two categories: *photorealistic* and *artistic*. Photorealistic style transfer [29, 32, 36, 56] aims at applying reference styles on scenes without hurting details and satisfying contradictory objectives. By contrast, artistic style transfer [4, 14, 20, 27, 28, 30, 35] captures style concepts from reference and modifies color distributions and

texture patterns of content images. However, it requires preparing numerous style images in advance, which limits practicality of style transfer. To tackle this issue, LDAST allows following textual descriptions to perform *artistic* style transfer and improves the accessibility of VFX design.

**Language-Based Image Editing.** The general task of LDAST is language-based image editing (LBIE), which also uses language to edit input images. With rule-based instructions and predefined semantic labels, they [7, 25] first carry out LBIE but under limited practicality. Inspired by text-to-image generation [40, 54, 58], previous works [5, 9,10,11, 26, 31, 33, 52] perform LBIE by conditional GAN, which modifies the properties of objects in the image. In contrast, LDAST aims at preserving the scene structure from the content image and performing stylization guided by the style instruction.

**CLIP-Guided Optimization.** Recently, based on the powerful visual-linguistic connection of CLIP [44], CLIP-guided image synthesis [34, 39] has shown exciting results. StyleCLIP [37] and NADA [12] tweak the latent code of a pre-trained StyleGAN [23] for image editing. Since heavily relying on a pre-trained generator, both are confined to the training domain, and the results can only present limited stylization. CLIPstyler [24] updates the style transfer network for target style patterns from the CLIP alignment.

Though supporting arbitrary content images, CLIPstyler still requires hundreds of iterations and takes lots of time with considerable GPU memory, suffering from the efficiency and practicality overhead. Moreover, our experiments show that CLIP poorly captures detailed style patterns from instructions, which is intractable to perform explicit LDAST.

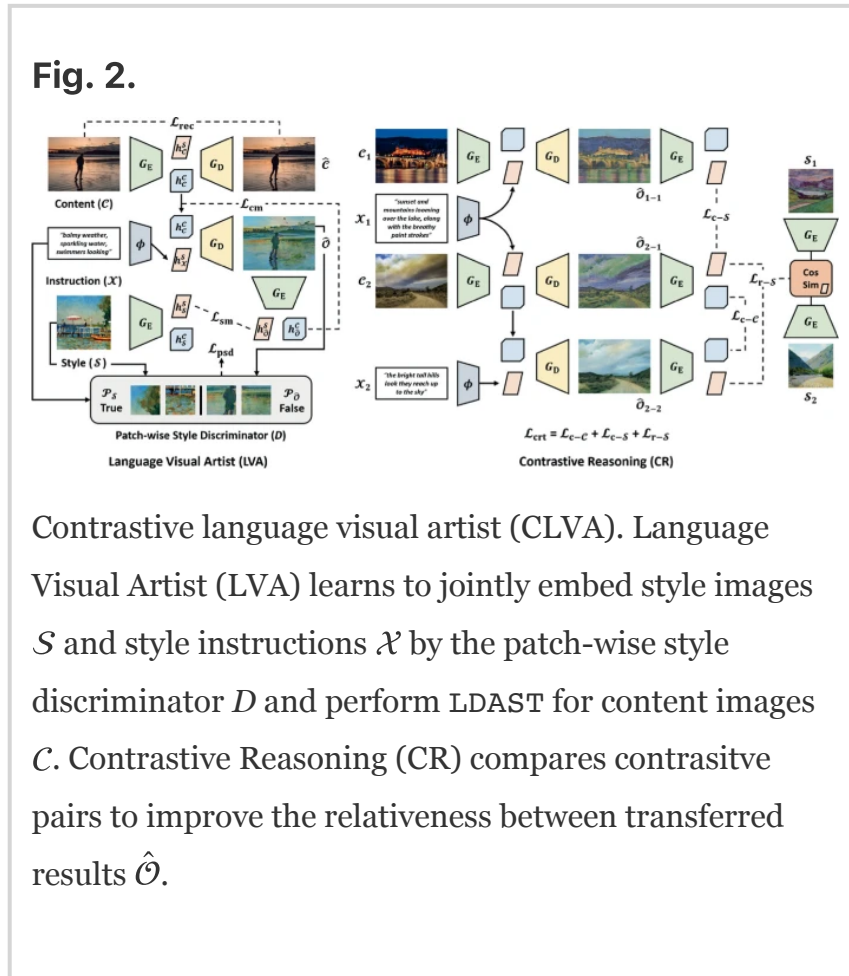
### 3 Language-Driven Artistic Style Transfer

---

#### 3.1 Overview of CLVA

We introduce language-driven artistic style transfer (LDAST) to manipulate the style of a content image  $C$ , guided by a style instruction  $\mathcal{X}$ , as illustrated in Fig. 1. For training, we have pairs of style images  $S$  with style instructions  $\mathcal{X}$  to learn the mutual correlation. During testing, only  $\mathcal{X}$  are provided for LDAST to carry out artistic style transfer purely relied on language. We present contrastive language visual artist (CLVA) in Fig. 2. Language visual artist (LVA) extracts content structures from  $C$  and visual patterns from  $\mathcal{X}$  to perform LDAST. LVA adopts the patch-wise style discriminator  $D$  to connect extracted visual semantics to patches of paired style image ( $\mathcal{P}_S$  in Fig. 2). Contrastive reasoning (CR) allows comparing contrastive pairs  $C_1-\mathcal{X}_1$ ,  $C_2-\mathcal{X}_1$ , and  $C_2-\mathcal{X}_2$  of content image and style instruction. In this way, it should present consistent content structures from the same content image  $C_2$  or

analogous style patterns from related style images  $S_1$  and  $S_2$ , despite using different style instructions.



Contrastive language visual artist (CLVA). Language Visual Artist (LVA) learns to jointly embed style images  $S$  and style instructions  $\mathcal{X}$  by the patch-wise style discriminator  $D$  and perform LDAST for content images  $C$ . Contrastive Reasoning (CR) compares contrastive pairs to improve the relativeness between transferred results  $\hat{O}$ .

### 3.2 Language Visual Artist (LVA)

To tackle LDAST, language visual artist (LVA) first adopts visual encoder  $G_E$  to extract the content feature  $h^C$  and the style feature  $h^S$  for an image. Text encoder  $\phi$  also extracts the style instruction feature  $h^S_{\mathcal{X}}$  from an instruction.  $h^C$  is a spatial tensor containing the content structure feature, and  $h^S$  represents the global style pattern.  $S^S_{\mathcal{X}}$  embeds into the same space of  $h^S$  to reflect the extracted visual semantic. Then, visual decoder  $G_D$  produces

transferred results  $\hat{\mathcal{O}}$  from  $h_C^C$  and  $h_X^S$ , which performs style transfer by style instructions:

$$\begin{aligned} h_C^C, h_C^S &= G_E(C), \quad h_X^S = \phi(\mathcal{X}), \\ \hat{\mathcal{O}} &= G_D(h_C^C, h_X^S). \end{aligned} \quad (1)$$

In particular,  $G_D$  applies self-attention [35, 57] to fuse  $h_C^C$  and  $h_X^S$  over the global spatial dimension.

There are two goals to train LVA for LDAST: (i) preserving *content structures* from content images; (ii) presenting *style patterns* correlated with visual semantics of style instructions.

**Structure Reconstruction.** To preserve content structures, we consider that visual decoder  $G_D$  should be able to reconstruct input content images using extracted content features  $h_C^C$  and style features  $h_C^S$  from visual encoder  $G_E$ :

$$\begin{aligned} \hat{C} &= G_D(h_C^C, h_C^S), \\ \mathcal{L}_{\text{rec}} &= \|\hat{C} - C\|_2, \end{aligned} \quad (2)$$

where the reconstruction loss  $\mathcal{L}_{\text{rec}}$  is computed as the mean L2 difference between reconstructed content images  $\hat{C}$  and input content images  $C$ .

**Patch-Wise Style Discriminator ( $D$ ).** Regarding style patterns, results  $\hat{\mathcal{O}}$  guided by style instructions

$\mathcal{X}$  are expected to present analogously to reference style images  $S$ . To address the connection between linguistic from  $\mathcal{X}$  and visual semantics from  $S$ , we introduce the patch-wise style discriminator  $D$ .

Inspired by texture synthesis [15, 53], images with analogous patch patterns should appear perceptually similar texture patterns.  $D$  tries to recognize the correspondence between an image patch  $\mathcal{P}$  and a style instruction  $\mathcal{X}$ :

$$\begin{aligned} \mathcal{P}_{\hat{\theta}}, \mathcal{P}_S &= \text{Crop}(\hat{\theta}), \text{Crop}(S), \\ \mathcal{L}_{\text{psd}} &= \log(1 - D(\mathcal{P}_{\hat{\theta}}, \mathcal{X})), \\ \mathcal{L}_D &= \log(1 - D(\mathcal{P}_{\hat{\theta}}, \mathcal{X})) + \log(D(\mathcal{P}_S, \mathcal{X})), \end{aligned} \tag{3}$$

where `Crop` is to randomly crop an image into patches. The patch-wise style loss  $\mathcal{L}_{\text{psd}}$  aims at generating transferred results that are correlated with  $\mathcal{X}$ . Contrarily, by the discriminator loss  $\mathcal{L}_D$ ,  $D$  learns to distinguish that a patch  $\mathcal{P}$  is from style images ( $\mathcal{P}_S$ ) or transferred results ( $\mathcal{P}_{\hat{\theta}}$ ). This adversarial loss [17, 41] encourages that transferred results from style instructions are presented similarly with style images, which jointly embeds the extracted visual semantics.

**Content Matching and Style Matching.** To further enhance the alignment with inputs, inspired by cycle consistency [38, 51, 55, 60], we consider the content matching loss  $\mathcal{L}_{\text{cm}}$  and the style matching

loss  $\mathcal{L}_{sm}$  of transferred results  $\hat{\mathcal{O}}$ . We adopt  $G_E$  again to extract content features  $h_{\hat{\mathcal{O}}}^C$  and style features  $h_{\hat{\mathcal{O}}}^S$  for  $\hat{\mathcal{O}}$ , where  $h_{\hat{\mathcal{O}}}^C$  and  $h_{\hat{\mathcal{O}}}^S$  should correlate with  $h_C^C$  from  $\mathcal{C}$  and  $h_S^S$  from  $\mathcal{S}$ :

$$(h_{\hat{\mathcal{O}}}^C, h_{\hat{\mathcal{O}}}^S), (h_S^C, h_S^S) = G_E(\hat{\mathcal{O}}), G_E(\mathcal{S}),$$

$$\mathcal{L}_{cm}, \mathcal{L}_{sm} = \|h_{\hat{\mathcal{O}}}^C - h_C^C\|_2, \|h_{\hat{\mathcal{O}}}^S - h_S^S\|_2.$$

(4)

Therefore, transferred results are required to align with content structures and style patterns from inputs, which meets the goal of LDAST.

### 3.3 Contrastive Reasoning (CR)

The content image should transfer to various styles while preserving the same structure. Related style instructions can apply analogous style patterns to arbitrary content images. As shown in Fig. 2, contrastive reasoning (CR) compares content structures or style patterns from transferred results of contrastive pairs. The contrastive pair consists of two different content images  $\mathcal{C}_1$  and  $\mathcal{C}_2$  with two reference styles  $\{\mathcal{S}_1, \mathcal{X}_1\}$  and  $\{\mathcal{S}_2, \mathcal{X}_2\}$ . We follow the LVA inference to acquire cross results for pairs of content images and style instructions:

$$\begin{aligned}
(h_{C_1}^C, h_{C_1}^S), (h_{C_2}^C, h_{C_2}^S) &= G_E(C_1), G_E(C_2), \\
h_{\mathcal{X}_1}^S, h_{\mathcal{X}_2}^S &= \phi(\mathcal{X}_1), \phi(\mathcal{X}_2), \\
\hat{\mathcal{O}}_{C_1-\mathcal{X}_1}, \hat{\mathcal{O}}_{C_1-\mathcal{X}_2} &= G_D(h_{C_1}^C, h_{\mathcal{X}_1}^S), G_D(h_{C_1}^C, h_{\mathcal{X}_2}^S), \\
\hat{\mathcal{O}}_{C_2-\mathcal{X}_1}, \hat{\mathcal{O}}_{C_2-\mathcal{X}_2} &= G_D(h_{C_2}^C, h_{\mathcal{X}_1}^S), G_D(h_{C_2}^C, h_{\mathcal{X}_2}^S).
\end{aligned}$$

(5)

**Consistent Matching.** Transferred results should present similar content structures ( $\hat{\mathcal{O}}_{C_2-\mathcal{X}_1}$  and  $\hat{\mathcal{O}}_{C_2-\mathcal{X}_2}$ ) or analogous style patterns ( $\hat{\mathcal{O}}_{C_1-\mathcal{X}_1}$  and  $\hat{\mathcal{O}}_{C_2-\mathcal{X}_1}$ ) if using the same content image ( $C_2$ ) or the same style instruction ( $\mathcal{X}_1$ ):

$$\begin{aligned}
h_{\hat{\mathcal{O}}_{C_i-\mathcal{X}_j}}^C &= G_E(\hat{\mathcal{O}}_{C_i-\mathcal{X}_j}), \\
\mathcal{L}_{C-C} &= \|h_{\hat{\mathcal{O}}_{C_1-\mathcal{X}_1}}^C - h_{\hat{\mathcal{O}}_{C_1-\mathcal{X}_2}}^C\|_2 + \|h_{\hat{\mathcal{O}}_{C_2-\mathcal{X}_1}}^C - h_{\hat{\mathcal{O}}_{C_2-\mathcal{X}_2}}^C\|_2 \\
\mathcal{L}_{C-S} &= \|h_{\hat{\mathcal{O}}_{C_1-\mathcal{X}_1}}^S - h_{\hat{\mathcal{S}}_{2-1}}^S\|_2 + \|h_{\hat{\mathcal{O}}_{C_1-\mathcal{X}_2}}^S - h_{\hat{\mathcal{O}}_{C_2-\mathcal{X}_2}}^S\|_2
\end{aligned}$$

(6)

where *consistent matching* of content structure  $\mathcal{L}_{C-C}$  or style pattern  $\mathcal{L}_{C-S}$  is aligned by content features or style features, extracted by  $G_E$ .

---

**Algorithm 1.** Training Process of Language Visual Artist (LVA)

---

```

1:  $G_E, G_D$ : Visual Encoder, Visual Decoder
2:  $\phi$ : Text Encoder
3:  $D$ : Patch-wise Style Discriminator
4: while TRAIN_VLA do
5:    $\mathcal{C}, \{\mathcal{S}, \mathcal{X}\} \leftarrow$  Sampled content/style
6:
7:    $h_{\mathcal{C}}^C, h_{\mathcal{C}}^S \leftarrow G_E(\mathcal{C})$     $\hat{\mathcal{C}} \leftarrow G_D(h_{\mathcal{C}}^C, h_{\mathcal{C}}^S)$ 
8:    $\mathcal{L}_{rec} \leftarrow$  Reconstruction loss
9:    $h_{\mathcal{X}}^S \leftarrow \phi(\mathcal{X})$     $\hat{\mathcal{O}} \leftarrow G_D(h_{\mathcal{C}}^C, h_{\mathcal{X}}^S)$  ▷ Eq. 2
10:   $\mathcal{P}_S, \mathcal{P}_{\hat{\mathcal{O}}} \leftarrow$  Crop( $\mathcal{S}$ ), Crop( $\hat{\mathcal{O}}$ )
11:   $\mathcal{L}_{psd} \leftarrow$  Patch-wise style loss ▷ Eq. 3
12:   $(h_{\hat{\mathcal{O}}}^C, h_{\hat{\mathcal{O}}}^S), (h_{\hat{\mathcal{S}}}^C, h_{\hat{\mathcal{S}}}^S) \leftarrow G_E(\hat{\mathcal{O}}), G_E(\mathcal{S})$ 
13:   $\mathcal{L}_{cm} \leftarrow$  Content matching loss ▷ Eq. 4
14:   $\mathcal{L}_{sm} \leftarrow$  Style matching loss ▷ Eq. 4
15:
16:   $\mathcal{L}_G \leftarrow \mathcal{L}_{rec} + \mathcal{L}_{psd} + \mathcal{L}_{cm} + \mathcal{L}_{sm}$  ▷ Eq. 3
17:   $\mathcal{L}_D \leftarrow$  Discriminator loss for D
18:  Update  $G_E, G_D, \phi$  by minimizing  $\mathcal{L}_G$ 
19:  Update  $D$  by maximizing  $\mathcal{L}_D$ 
20: end while

```

---



**Relative Matching.** Apart from consistent matching, distinct style instructions, which imply corresponding visual semantics, should still present relative patterns. For example, we can only discover “*reach up to the sky*” literally from  $\mathcal{X}_2$ . If comparing reference style images  $\mathcal{S}_1$  and  $\mathcal{S}_2$ , we can perceive the sharing of a similar style pattern and link the visual concept of “*bright tall hills*” in  $\mathcal{X}_2$  to “*mountains looming over the lake*” in  $\mathcal{X}_1$ . We define *relative matching*  $\mathcal{L}_{r-S}$  with the cosine similarity (CosSim) between reference style images:

$$\begin{aligned}
 (h_{\mathcal{S}_i}^C, h_{\mathcal{S}_i}^S) &= G_E(\mathcal{S}_i), \\
 r &= \text{CosSim}(h_{\mathcal{S}_1}^S, h_{\mathcal{S}_2}^S), \\
 \mathcal{L}_{r-S} &= (\|h_{\hat{\theta}_{c_1-\mathcal{X}_1}}^S - h_{\hat{\theta}_{c_1-\mathcal{X}_2}}^S\|_2 + \\
 &\quad \|h_{\hat{\theta}_{c_2-\mathcal{X}_1}}^S - h_{\hat{\theta}_{c_2-\mathcal{X}_2}}^S\|_2) \cdot r.
 \end{aligned}
 \tag{7}$$

When style images are related, it has to align style features to certain extent even if paired style instructions are different. Otherwise,  $\mathcal{L}_{r-S}$  will be close to 0 and ignore this unrelated style pair. The overall contrastive reasoning loss  $\mathcal{L}_{\text{ctr}}$  considers both consistent matching and relative matching:

$$\mathcal{L}_{\text{ctr}} = \mathcal{L}_{c-C} + \mathcal{L}_{c-S} + \mathcal{L}_{r-S}.
 \tag{8}$$

### 3.4 Learning of CLVA

For each epoch of CLVA training, we first train with the LVA process and then CR. As Algorithm 1, we consider reconstruction loss  $\mathcal{L}_{\text{rec}}$  to preserve content structure and patch-wise style loss  $\mathcal{L}_{\text{psd}}$  between style instruction and visual pattern of transferred results. Both content matching loss  $\mathcal{L}_{\text{cm}}$  and style matching loss  $\mathcal{L}_{\text{sm}}$  enhance the matching with the inputs. Simultaneously, we update  $D$  by maximizing discriminator loss  $\mathcal{L}_D$  to distinguish between true patches  $\mathcal{P}_S$  or false patches  $\mathcal{P}_{\hat{\theta}}$ , concerning style instructions. During CR, contrastive pairs of content images and style instructions are randomly sampled, and the transferred results are across produced. We further update by minimizing contrastive reasoning loss  $\mathcal{L}_{\text{ctr}}$  to allow considering content consistency and mutual style relativeness. The overall optimization of CLVA is summarized as:

$$\begin{aligned} \mathcal{L}_G &= \mathcal{L}_{\text{rec}} + \mathcal{L}_{\text{psd}} + \mathcal{L}_{\text{cm}} + \mathcal{L}_{\text{sm}}, \\ \min_{G, \phi} \max_D & \mathcal{L}_G + \mathcal{L}_D + \mathcal{L}_{\text{ctr}}. \end{aligned}$$

(9)

## 4 Experiments

---

### 4.1 Experimental Setup

**Dataset.** To evaluate our CLVA, we consider DTD<sup>2</sup> [8] and ArtEmis [1] as reference style instructions. DTD<sup>2</sup> contains 5K texture images with its natural descriptions for visual attributes such as

colors and texture patterns. ArtEmis provides 80K artworks from WikiArt<sup>1</sup> with annotations of visual contents and emotional effects as human style feelings. We also collect 15K wallpapers from WallpapersCraft<sup>2</sup>, which presents diverse scenes as content images. Each content image is resized to  $256 \times 192$  in our experiment. We randomly sample 100 unseen content images and 100 testing reference styles to evaluate the generalizability of LDAST. Note that both style images and style instructions appear for training, but only style instructions are accessible during testing.

**Evaluation Metrics.** To support large-scale evaluation, we treat transferred results directly from style images as semi-groundtruth (Semi-GT) [2, 3, 42] by the SOTA style transfer AdaAttn [30]. We apply the following metrics:

- **SSIM** [48] compares images in the luminance, contrast, and structure aspects. A higher SSIM has a higher structural similarity;
- **Percept** [22] computes from the gram matrix of visual features. A lower Percept loss shows that two images share a similar style pattern;
- **FAD** [18] is computed by the mean L2 distance of the activations from the InceptionV3 [46] feature. As a distance metric, a lower FAD represents that LDAST results and Semi-GT are more relevant.

Note that we consider SSIM and FAD to compare with Semi-GT and calculate Percept loss directly with reference style images. Apart from visual similarity, we consider the correlation between style instructions and LDAST results:

- **VLS** [49] calculates the cosine similarity between each other from CLIP [44].

Since each metric has different deficiencies, we also conduct a comprehensive human evaluation from aspects of content, instruction, and style matching. We randomly sample 75 LDAST results and adopt MTurk<sup>3</sup> to rank over all methods. We also hire 3 MTurkers for each task to avoid the potential ranking bias.

**Baselines.** We conduct baselines for LDAST from various aspects:

- **Style Transfer:** We consider previous artistic style transfer methods NST [14], WCT [28], Adaln [20], SANet [35], and LST [27] that support arbitrary content images. We use the same style (instruction and image) encoding from our CLVA as style features and follow their own training process to perform LDAST upon them. Due to the space issue, we only show the comparison with more recent SANet and LST. Please refer to Appendix for the complete results.

- **Language-based Image Editing:** We adopt ManiGAN [26] with affine combination module (ACM) as the general language-based editing baseline, where it modifies the content image by the style instruction. We treat normal style transferred results as groundtruth for ManiGAN to learn from.
- **CLIP-based Optimization:** StyleCLIP [37], NADA [12], and CLIPstyler [24] manipulate the content image based on the CLIP alignment of the guided instruction. Since StyleCLIP and NADA are restricted by the pre-trained generator, we compare them with the training domains of car and church. Differently, CLIPstyler can carry out arbitrary content images for LDAST.

## 4.2 Quantitative Results

**Instruction with Visual Attributes.** Table 1 illustrates the comparison of LDAST with baselines on DTD<sup>2</sup>. As regards automatic metrics, CLVA preserves content structures (highest 36.65 SSIM) and stylizes with related visual attributes to style images (lowest 0.2033 Percept loss). Furthermore, CLVA brings out the highest overall similarity as Semi-GT (lowest 0.1493 FAD). Since CLIPstyler directly optimizes by CLIP [44], it makes the highest VLS. Through the patch-wise discriminator, our CLVA can still produce style patterns correlated to given

instructions (competitive 24.00 VLS) even without the pre-trained CLIP.

**Table 1. Testing results of  $\text{LDAST}$  using visual attribute instructions on  $\text{DTD}^2$ .**

**Fig. 3.**



Visualized comparison using visual attribute instructions on  $\text{DTD}^2$ .

The human evaluation investigates the matching between transferred results with content images (Content), style instructions (Instruction), style images (Style), and Semi-GT (Semi-GT). In particular, content and instruction matching are the two most crucial, which concern the goal of  $\text{LDAST}$ : *content structure preservation* and *style pattern presentation*; style image and semi-gt matching are provided for different comparing targets from a human aspect. The results are calculated by the mean ranking score (from 1 to 5, the higher is better) of each method. In general, MTurkers indicate that our CLVA has an apparent advantage in preserving

content structures (highest 3.852 Content) and presenting aligned style patterns (highest 3.742 Instruction). Though with the aid of CLIP, CLIPstyler is still behind CLVA (-0.4 Instruction), with an even higher gap in style image matching (-0.5 Style). Contributed by contrastive reasoning that compares the mutual relativeness between pairs of contents and instructions, CLVA can stylize with the captured visual attributes. We adopt Pearson correlation and investigate the coefficients between automatic metrics and human evaluation as 77.2 (FAD  $\rightarrow$  Instruction), 84.5 (FAD  $\rightarrow$  Semi-GT), 81.3 (VLS  $\rightarrow$  Instruction), and 77.8 (VLS  $\rightarrow$  Semi-GT). This high correlation indicates that our metric design is adequate for evaluating large-scale LDATA experiments. The even higher 88.2 correlation (Instruction  $\rightarrow$  Semi-GT) between instruction and Semi-GT matching in human evaluation further supports the usage of Semi-GT.

From the aspect of visualized comparison in Fig. 3, previous SANet and LST only produce repetitive and disorder textures in their transferred results.

ManiGAN modifies the style directly over pixels, suffering from blurring objects; this deficiency can also be found in Table 1 (lower SSIM and lower Content matching). CLIPstyler is sometimes misguided by CLIP, making irrelevant patterns, such as the bright white background in the third case. Contrary to baselines, CLVA extracts a more detailed

style from different kinds of guidance ("*brown metallic*" in the first row and "*stringy hairy*" in the third case), leading to superior LDAST results that correspond to style instructions.

**Table 2. Testing results of LDAST using emotional effect instructions on ArtEmis.**

**Fig. 4.**



Visualized comparison using emotional effect instructions on ArtEmis.

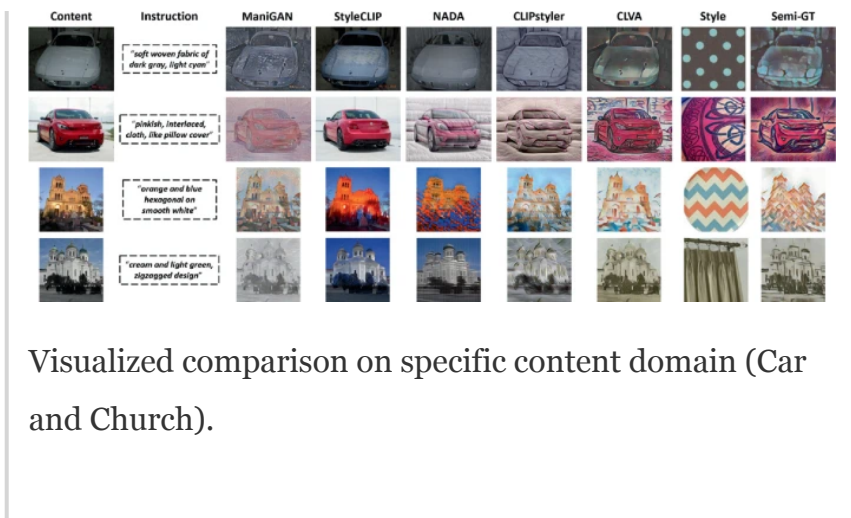
**Instruction with Emotional Effects.** Unlike visual attributes, emotional effect instructions are more challenging as connecting to visual semantics of described objects or style patterns from human feelings. For example, "*yellowish and green*" from "*sunset and mountains*" or "*scaring charcoal grey*" from "*nightmare*". We consider this human style feeling on ArtEmis [1], where the model has to express the latent visual concepts of emotional effect instructions. CLVA performs with more balance (both second-highest SSIM and second-



lowest Percept) from Table 2, especially the lowest 0.1418 FAD, making the most similar transferred results to Semi-GT. Though CLIPstyler [24] achieves higher VLS by optimizing over CLIP, from human aspects, CLVA can preserve more concrete contents and present more correlated style patterns (higher 3.357 content and 3.586 instruction matching). The visualized comparison in Fig. 4 illustrates that previous SANet [35] and LST [27] contain unsmooth and fragmentary patterns with blurring contents. Without a style transformation process, ManiGAN [26] modifies with only monotonous colors. CLIPstyler is failed to capture human style feelings well, suffering from weird and unpleasant results. Different from them, our CLVA learns the visual semantic during contrastive reasoning by comparing mutual relativeness between literal instructions and style images, leading to a more colorful and corresponding stylization as human emotion. More surprisingly, despite not instructed literally, CLVA perceives "*side of the water*" and reveals the latent yet correlated "*grassland*" precisely in the third row.

**Table 3. Testing results of LDAST on specific content domain (Car and Church).**

**Fig. 5.**



**Specific Content Domain.** To compare with StyleCLIP [37] and NADA [12] that are restricted by the pre-trained generator, we evaluate LDATA on the specific content domain. We consider the same domain images in StyleGAN2 [23] and visual attribute instructions on DTD<sup>2</sup>. Table 3 indicates the numerical comparison on Car and Church. Our CLVA still produces superior results and is the most admirable by human. Since StyleCLIP and NADA rely on StyleGAN, they can only preserve content (highest 3.459 Content by StyleCLIP) but with limited stylization (lower Instruction and Style). Similar observations can be found in Fig. 5, where StyleCLIP shows almost no modification for the second car. They can neither deal with the background; NADA even destroys the scene in the third row. In contrast to CLIPstyer [24] that only contains abstractive and obscure styles, CLVA presents the detailed “read interplaced cloth” behind the car and the color “cream” precisely on the surface of the church.

**Table 4. Ablation study of CLVA using visual attribute instructions on DTD<sup>2</sup>.**

**Table 5. Instruction-to-style retrieval on DTD<sup>2</sup> and ArtEmis.**

**Table 6. Human comparison between CLVA and CLIPstyle with fine-tuned CLIP on DTD<sup>2</sup>.**

**Fig. 6.**



Visualization examples of instruction-to-style retrieval by CLIP and CLVA.

### 4.3 Ablation Study

We conduct an ablation study of each component effect on DTD<sup>2</sup> in Table 4. At row (a), with the reconstruction  $\mathcal{L}_{\text{rec}}$  and the patch-wise style  $\mathcal{L}_{\text{psd}}$ , CLVA achieves feasible LDAST results by concrete structures and extracted style semantics. Row (b)–(d) shows the strength of content matching  $\mathcal{L}_{\text{cm}}$  and style matching  $\mathcal{L}_{\text{sm}}$ . In particular, content matching

helps the structure similarity to content images (higher 36.05 SSIM). Style matching aims at analogous visual patterns to style images, which leads to better stylization quality (lower 0.2049 Percept and higher 23.69 VLS). If considering altogether, it can benefit and strike a balance between both. Finally, contrastive reasoning  $\mathcal{L}_{ctr}$  further enables CLVA to consider contrastive pairs, making a comprehensive improvement at row (e).

**Table 7. Time and GPU cost when performing LDAST on TITAN X with content image size  $256 \times 192$ . \* means this method can only run one input at a time.**

**Fig. 7.**



Style interpolation results of LDAST over instructions.

**Why CLVA is Better than CLIP-Based?** Despite no CLIP optimized, CLVA demonstrates superior results on LDAST with all aspects of automatic metrics and human evaluation. To investigate it, we conduct instruction-to-style retrieval based on the similarity between features of style instructions and style

images. Table 5 shows that our learned CLVA performs higher Recall@k on both DTD<sup>2</sup> and ArtEmis, leading to a better instruction-style alignment than the used CLIP. The visualization in Fig. 6 also indicates the flaw of CLIP on detailed style patterns. For example, in the first row, CLIP only presents either “*bright color*” or “*town*” in the retrieval results. In contrast, CLVA can capture both and present more related LDAST to “*happy place to live*”. From Table 6, even CLIP has been fine-tuned ahead; our CLVA still produces preferable LDAST results from all human aspects of content, instruction, and style matching. This observation supports that contrastive reasoning, which considers contrastive pairs of content images and style instructions, is required to benefit from mutual relateness.

Apart from transfer quality, CLVA also holds a higher efficiency than CLIP-based methods. Table 7 illustrates the time and GPU cost on a single TITAN X (12GB) with content image size 256 × 192. All CLIP-based methods take more than 30 s for only one pair of content images and style instructions. Instead of numerous iterations to align with CLIP, we extract style semantics and carry out LDAST in one shot, taking merely 0.03 s for one input. Without updating the model during inference, CLVA supports parallelization and can accomplish 50 pairs in half a second. Besides, as a lightweight style transfer

network, CLVA requires the least GPU memory for LDAST. In summary, our CLVA surpasses those CLIP-based methods on both quality and efficiency because of the detailed style deficiency and the required optimizing iteration from CLIP.

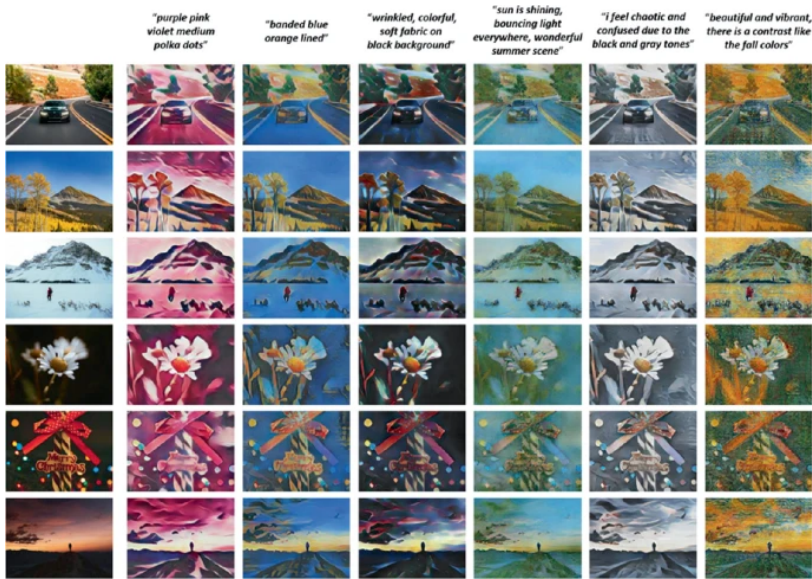
**Qualitative Results.** As shown in Fig. [7](#), we investigate the linear interpolation of extracted style patterns by CLVA. Considering style features  $h_{\mathcal{X}_1}^S$  and  $h_{\mathcal{X}_2}^S$  of instructions  $\mathcal{X}_1$  and  $\mathcal{X}_2$ , the interpolated  $h_p^S$  should be:

$$h_p^S = (1 - \alpha)h_{\mathcal{X}_1}^S + \alpha h_{\mathcal{X}_2}^S, \quad (10)$$

where  $\alpha$  is the style ratio between the two. Figure [7](#) presents a smooth transformation from one style instruction to another. By training on DTD<sup>2</sup> and ArtEmis altogether, CLVA even performs interpolated stylization by both visual attribute and emotional effect instructions in the third row. Figure [8](#) illustrates diverse LDAST results by our CLVA. Since CLVA supports arbitrary content images, we can also modify the style detail for high-resolution inputs in Fig. [9](#).

**Fig. 8.**





CLVA results on diverse pairs of content images and style instructions.

Fig. 9.



High-resolution (1920 × 1080) LDAST results by CLVA with upper right: *“the lonely world makes me feel scared and nostalgic how sky and sea merge together”*; lower left: *“the snow and lights in the shop windows looks like a winter scene”*; lower right: *“ink painting, black dotted line, whiteboard”*.

## 5 Conclusion

---

We introduce language-driven artistic style transfer (LDAST) to do stylization for a content image by a style instruction. We propose contrastive language visual artist (CLVA) that adopts the patch-wise style discriminator and contrastive reasoning to jointly learn between style images and style instructions. We demonstrate that CLVA can express various style patterns of visual attributes as well as emotional effects and perform LDAST efficiently. CLVA also outperforms baselines on both automatic metrics and human evaluation. We believe that LDAST can make visual applications like image/video effect more controllable for humans.

## Notes

---

1. WikiArt: <https://www.wikiart.org>.
2. WallpapersCraft: <https://wallpaperscraft.com/>.
3. Amazon Mechanical Turk: <https://www.mturk.com>.

## References

---

1. Achlioptas, P., Ovsjanikov, M., Haydarov, K., Elhoseiny, M., Guibas, L.: ArtEmis: affective language for visual art. In: CVPR (2021)
-



2. Al-Sarraf, A., Shin, B.-S., Xu, Z., Klette, R.:  
Ground truth and performance evaluation of lane  
border detection. In: Chmielewski, L.J., Kozera,  
R., Shin, B.-S., Wojciechowski, K. (eds.) ICCVG  
2014. LNCS, vol. 8671, pp. 66–74. Springer,  
Cham (2014). [https://doi.org/10.1007/978-3-319-11331-9\\_9](https://doi.org/10.1007/978-3-319-11331-9_9)

---

3. Borkar, A., Hayes, M., Smith, M.T.: An efficient  
method to generate ground truth for evaluating  
lane detection systems. In: ICASSP (2010)

---

4. Chen, H., et al.: DualAST: dual style-learning  
networks for artistic style transfer. In: CVPR  
(2021)

---

5. Chen, J., Shen, Y., Gao, J., Liu, J., Liu, X.:  
Language-based image editing with recurrent  
attentive models. In: CVPR (2018)

---

6. Chen, Y.L., Hsu, C.T.: Towards deep style  
transfer: a content-aware perspective. In: BMVC  
(2016)

---

7. Cheng, M.M., et al.: ImageSpirit: verbal guided  
image parsing. In: ACM Transactions on Graphics  
(2013)

---

8. Cimpoi, M., Maji, S., Kokkinos, I., Mohamed, S., Vedaldi, A.: Describing textures in the wild. In: CVPR (2014)

---
9. El-Nouby, A., et al.: Tell, draw, and repeat: generating and modifying images based on continual linguistic instruction. In: ICCV (2019)

---
10. Fu, T.J., Wang, X.E., Grafton, S., Eckstein, M., Wang, W.Y.: SSCR: iterative language-based image editing via self-supervised counterfactual reasoning. In: EMNLP (2020)

---
11. Fu, T.J., Wang, X.E., Grafton, S., Eckstein, M., Wang, W.Y.: M3L: language-based video editing via multi-modal multi-level transformer. In: CVPR (2022)

---
12. Gal, R., Patashnik, O., Maron, H., Chechik, G., Cohen-Or, D.: StyleGAN-NADA: CLIP-guided domain adaptation of image generators. [arXiv:2108.00946](https://arxiv.org/abs/2108.00946) (2021)

---
13. Gao, C., Gu, D., Zhang, F., Yu, Y.: ReCoNet: real-time coherent video style transfer network. [arXiv:1807.01197](https://arxiv.org/abs/1807.01197) (2018)

---
14. Gatys, L.A., Ecker, A.S., Bethge, M.: A neural algorithm of artistic style. [arXiv:1508.06576](https://arxiv.org/abs/1508.06576)

(2015)

---

15. Gatys, L.A., Ecker, A.S., Bethge, M.: Texture synthesis using convolutional neural networks. In: NeurIPS (2015)

---

16. Gatys, L.A., Ecker, A.S., Bethge, M., Hertzmann, A., Shechtman, E.: Controlling perceptual factors in neural style transfer. In: CVPR (2017)

---

17. Goodfellow, I.J., et al.: Generative adversarial networks. In: NeurIPS (2014)

---

18. Heusel, M., Ramsauer, H., Unterthiner, T., Nessler, B., Hochreiter, S.: GANs trained by a two time-scale update rule converge to a local nash equilibrium. In: NeurIPS (2017)

---

19. Huang, H., et al.: Real-time neural style transfer for videos. In: CVPR (2017)

---

20. Huang, X., Belongie, S.: Arbitrary style transfer in real-time with adaptive instance normalization. In: ICCV (2017)

---

21. Jing, Y., Yang, Y., Feng, Z., Ye, J., Yu, Y., Song, M.: Neural style transfer: a review. [arXiv:1705.04058](https://arxiv.org/abs/1705.04058) (2017)

---

22. Johnson, J., Alahi, A., Fei-Fei, L.: Perceptual losses for real-time style transfer and super-resolution. In: ECCV (2016)

---
23. Karras, T., Laine, S., Aittala, M., Hellsten, J., Lehtinen, J., Aila, T.: Analyzing and improving the image quality of StyleGAN. In: CVPR (2020)

---
24. Kwon, G., Ye, J.C.: CLIPstyler: image style transfer with a single text condition. In: CVPR (2022)

---
25. Laput, G., et al.: PixelTone: a multimodal interface for image editing. In: CHI (2013)

---
26. Li, B., Qi, X., Lukasiewicz, T., Torr, P.H.S.: ManiGAN: text-guided image manipulation. In: CVPR (2020)

---
27. Li, X., Liu, S., Kautz, J., Yang, M.H.: Learning linear transformations for fast arbitrary style transfer. In: CVPR (2019)

---
28. Li, Y., Fang, C., Yang, J., Wang, Z., Lu, X., Yang, M.H.: Universal style transfer via feature transforms. In: NeurIPS (2017)

---
29. Li, Y., Liu, M.Y., Li, X., Yang, M.H., Kautz, J.: A

closed-form solution to photorealistic image stylization. In: ECCV (2018)

---

30. Liu, S., et al.: AdaAttN: revisit attention mechanism in arbitrary neural style transfer. In: ICCV (2021)

---

31. Liu, X., et al.: Open-edit: open-domain image manipulation with open-vocabulary instructions. In: Vedaldi, A., Bischof, H., Brox, T., Frahm, J.-M. (eds.) ECCV 2020. LNCS, vol. 12356, pp. 89–106. Springer, Cham (2020).  
[https://doi.org/10.1007/978-3-030-58621-8\\_6](https://doi.org/10.1007/978-3-030-58621-8_6)

---

32. Luan, F., Paris, S., Shechtman, E., Bala, K.: Deep photo style transfer. In: CVPR (2017)

---

33. Nam, S., Kim, Y., Kim, S.J.: Text-adaptive generative adversarial networks: manipulating images with natural language. In: NeurIPS (2018)

---

34. Nichol, A., et al.: GLIDE: towards photorealistic image generation and editing with text-guided diffusion models. In: [arXiv:2112.10741](https://arxiv.org/abs/2112.10741) (2021)

---

35. Park, D.Y., Lee, K.H.: Arbitrary style transfer with style-attentional networks. In: CVPR

(2019)

---

36. Park, T., et al.: Swapping autoencoder for deep image manipulation. In: NeurIPS (2020)

---

37. Patashnik, O., Wu, Z., Shechtman, E., Cohen-Or, D., Lischinski, D.: StyleCLIP: text-driven manipulation of StyleGAN imagery. In: ICCV (2021)

---

38. Qiao, T., Zhang, J., Xu, D., Tao, D.: MirrorGAN: Learning Text-to-image Generation by Redescription. In: CVPR (2019)

---

39. Ramesh, A., et al.: Zero-shot text-to-image generation. In: [arXiv:2102.12092](https://arxiv.org/abs/2102.12092) (2021)

---

40. Reed, S., Akata, Z., Yan, X., Logeswaran, L., Schiele, B., Lee, H.: Generative adversarial text to image synthesis. In: ICML (2016)

---

41. Salehi, P., Chalechale, A., Taghizadeh, M.: Generative adversarial networks (GANs): an overview of theoretical model, evaluation metrics, and recent developments. [arXiv:2005.13178](https://arxiv.org/abs/2005.13178) (2020)

---

42. Salvo, R.D.: Large scale ground truth generation

for performance evaluation of computer vision methods. In: VIGTA (2013)

---

43. Sanakoyeu, A., Kotovenko, D., Lang, S., Ommer, B.: A style-aware content loss for real-time HD style transfer. In: ECCV (2018)

---

44. Shi, L., et al.: Contrastive visual-linguistic pretraining. [arXiv:2007.13135](https://arxiv.org/abs/2007.13135) (2020)

---

45. Somavarapu, N., Ma, C.Y., Kira, Z.: Frustratingly simple domain generalization via image stylization. [arXiv:2006.11207](https://arxiv.org/abs/2006.11207) (2020)

---

46. Szegedy, C., Vanhoucke, V., Ioffe, S., Shlens, J., Wojna, Z.: Rethinking the inception architecture for computer vision. In: CVPR (2016)

---

47. Wang, P., Li, Y., Vasconcelos, N.: Rethinking and improving the robustness of image style transfer. In: CVPR (2021)

---

48. Wang, Z., Bovik, A.C., Sheikh, H.R., Simoncel, E.P.: Image quality assessment: from error visibility to structural similarity. In: TIP (2004)

---

49. Wu, C., et al.: GODIVA: generating open-Domain videos from nAtural descriptions.

50. Wu, C., Timm, M., Maji, S.: Describing textures using natural language. In: Vedaldi, A., Bischof, H., Brox, T., Frahm, J.-M. (eds.) ECCV 2020. LNCS, vol. 12346, pp. 52–70. Springer, Cham (2020). [https://doi.org/10.1007/978-3-030-58452-8\\_4](https://doi.org/10.1007/978-3-030-58452-8_4)
- 
51. Wu, L., Wang, Y., Shao, L.: Cycle-consistent deep generative hashing for cross-modal retrieval. In: TIP (2018)
- 
52. Xia, W., Yang, Y., Xue, J.H., Wu, B.: TediGAN: text-guided diverse face image generation and manipulation. In: CVPR (2021)
- 
53. Xian, W., et al.: TextureGAN: controlling deep image synthesis with texture patches. In: CVPR (2018)
- 
54. Xu, T., et al.: AttnGAN: fine-grained text to image generation with attentional generative adversarial networks. In: CVPR (2018)
- 
55. Yi, Z., Zhang, H., Tan, P., Gong, M.: DualGAN: unsupervised dual learning for image-to-image translation. In: ICCV (2017)
-



56. Yoo, J., Uh, Y., Chun, S., Kang, B., Ha, J.W.: Photorealistic style transfer via wavelet transforms. In: ICCV (2019)
- 
57. Zhang, H., Goodfellow, I., Metaxas, D., Odena, A.: Self-attention generative adversarial networks. In: PMLR (2019)
- 
58. Zhang, H., et al.: StackGAN: text to photo-realistic image synthesis with stacked generative adversarial networks. In: ICCV (2017)
- 
59. Zhang, Z., Wang, Z., Lin, Z., Qi, H.: Image super-resolution by neural texture transfer. In: CVPR (2019)
- 
60. Zhu, J.Y., Park, T., Isola, P., Efros, A.A.: Unpaired image-to-image translation using cycle-consistent adversarial networks. In: ICCV (2017)
- 

## Acknowledgments

---

Research was sponsored by the U.S. Army Research Office and was accomplished under Contract Number W911NF-19-D-0001 for the Institute for Collaborative Biotechnologies. The views and

conclusions contained in this document are those of the authors and should not be interpreted as representing the official policies, either expressed or implied, of the U.S. Government. The U.S. Government is authorized to reproduce and distribute reprints for Government purposes notwithstanding any copyright notation herein.

## Author information

---

### Authors and Affiliations

**UC Santa Barbara, Santa Barbara, USA**

Tsu-Jui Fu & William Yang Wang

**UC Santa Cruz, Santa Cruz, USA**

Xin Eric Wang

Corresponding author

Correspondence to [Tsu-Jui Fu](#).

## Editor information

---

### Editors and Affiliations

**Tel Aviv University, Tel Aviv, Israel**

Shai Avidan

**University College London, London, UK**

Gabriel Brostow

**Google AI, Accra, Ghana**

Moustapha Cissé

**University of Catania, Catania, Italy**

Giovanni Maria Farinella

**Facebook (United States), Menlo Park, CA, USA**

Tal Hassner

## 1 Electronic supplementary material

---

Below is the link to the electronic supplementary material.

[Supplementary material 1 \(pdf 9698 KB\)](#)

## Rights and permissions

---

[Reprints and Permissions](#)

## Copyright information

---

© 2022 The Author(s), under exclusive license to Springer Nature Switzerland AG

## About this paper

---

### Cite this paper

Fu, T.J., Wang, X.E., Wang, W.Y. (2022). Language-Driven Artistic Style Transfer. In: Avidan, S., Brostow, G., Cissé, M., Farinella, G.M., Hassner, T. (eds) Computer Vision – ECCV 2022. ECCV 2022. Lecture Notes in Computer Science, vol 13696. Springer, Cham. [https://doi.org/10.1007/978-3-031-20059-5\\_41](https://doi.org/10.1007/978-3-031-20059-5_41)

[.RIS](#)↓ [.ENW](#)↓ [.BIB](#)↓

DOI

[https://doi.org/10.1007/978-3-031-20059-5\\_41](https://doi.org/10.1007/978-3-031-20059-5_41)

|                 |                |                   |
|-----------------|----------------|-------------------|
| Published       | Publisher Name | Print ISBN        |
| 29 October 2022 | Springer, Cham | 978-3-031-20058-8 |

|                   |                                  |
|-------------------|----------------------------------|
| Online ISBN       | eBook Packages                   |
| 978-3-031-20059-5 | <a href="#">Computer Science</a> |
|                   | <a href="#">Computer Science</a> |
|                   | <a href="#">(RO)</a>             |

## Share this paper

Anyone you share the following link with will be able to read this content:

[Get shareable link](#)

Provided by the Springer Nature SharedIt content-sharing initiative

Not logged in - 169.231.177.94

California Digital Library TA (3902410118) - California Digital Library (3000123641) - University of California, Santa Barbara (8200642617)

**SPRINGER NATURE**

© 2022 Springer Nature Switzerland AG. Part of [Springer Nature](#).

# Language-Driven Artistic Style Transfer

Tsu-Jui Fu<sup>†</sup>, Xin Eric Wang<sup>‡</sup>, William Yang Wang<sup>†</sup>

<sup>†</sup>UC Santa Barbara    <sup>‡</sup>UC Santa Cruz  
{tsu-juifu, william}@cs.ucsb.edu    xwang366@ucsc.edu

**Abstract.** Despite having promising results, style transfer, which requires preparing style images in advance, may result in lack of creativity and accessibility. Following human instruction, on the other hand, is the most natural way to perform artistic style transfer that can significantly improve controllability for visual effect applications. We introduce a new task—language-driven artistic style transfer (**LDAST**)—to manipulate the style of a content image, guided by a text. We propose contrastive language visual artist (CLVA) that learns to extract visual semantics from style instructions and accomplish LDAST by the patch-wise style discriminator. The discriminator considers the correlation between language and patches of style images or transferred results to jointly embed style instructions. CLVA further compares contrastive pairs of content images and style instructions to improve the mutual relativeness. The results from the same content image can preserve consistent content structures. Besides, they should present analogous style patterns from style instructions that contain similar visual semantics. The experiments show that our CLVA is effective and achieves superb transferred results on LDAST.

## 1 Introduction

Style transfer [14,28,20,35,27,21] adopts appearances and visual patterns from another reference style images to manipulate a content image. Artistic style transfer has a considerable application value for creative visual design, such as image stylization and video effect [45,59,13,19]. However, it requires preparing collections of style image in advance. It even needs to redraw new references first if there is no expected style images, which is impractical due to an additional overhead. In contrast, language is the most natural way for humans to communicate. If a system can follow textual descriptions and automatically perform style transfer, we can significantly improve accessibility and controllability.

In this paper, we introduce Language-driven Artistic Style Transfer (**LDAST**). As illustrated in Fig. 1, LDAST treats a content image and a text as the input, and the style transferred result is manipulated based on the style description. It should preserve the structure of the content yet simultaneously modifies the style pattern that corresponds to the instruction. LDAST is different from the general language-based image-editing (LBIE) [33,26,31,9] that aims at altering objects or properties of objects. The main challenge of LDAST is to extract visual semantics from language. Humans use not only explicit visual attributes but



**Fig. 1.** Language-driven Artistic Style Transfer (LDAST). LDAST performs style transfer for a content image  $\mathcal{C}$ , guided by the visual attribute (the lower row) or even the visual content and emotional effect (the upper row) from a style instruction  $\mathcal{X}$ .

also visual content or emotional effects to describe style feelings. For example, it requires connecting “*water, sketching, and painting*” or “*peaceful, feel content*” with their visual concepts and further carrying out correlated style transfer.

We present contrastive language visual artist (CLVA), including language visual artist (LVA) and contrastive reasoning (CR), to perform style transfer conditioning on guided texts. LVA preserves content structures from content images  $\mathcal{C}$  and extracts visual semantics from style instructions  $\mathcal{X}$ . LVA learns the latent style pattern based on the distinguishment between patches of style images or transferred results from the patch-wise style discriminator. Furthermore, CR boosts by comparing contrastive pairs where relative content images or style instructions should present similar content structures or style patterns.

To evaluate LDAST, we conduct experiments upon DTD<sup>2</sup> [50] and ArtEmis [1]. DTD<sup>2</sup> provides texture images with its colors or texture patterns in text. ArtEmis collects explanations of visual contents and emotional effects for artworks. We treat these annotations as style instructions for the challenging LDAST concerning visual attributes or human style feelings. The experiments show that our CLVA is effective for LDAST and achieves superb yet efficient transferred results on both automatic metrics and human evaluation. Our contributions are four-fold:

- We introduce LDAST that follows natural language for artistic style transfer;
- We present CLVA, which learns to extract explicit visual semantics from style instructions and provide sufficient style patterns for LDAST;
- We conduct the evaluation on DTD<sup>2</sup> and ArtEmis to consider diverse style instructions with visual attributes and emotional effects;
- Extensive experiments and qualitative examples demonstrate that our CLVA outperforms baselines regarding both effectiveness and efficiency.

## 2 Related Work

**Artistic Style Transfer.** Style transfer [14,21,47,6,16,22,43] redraws an image with a specific style. Since being a popular form of art, incorporating painting with digital design can produce attractive visual effect (VFX). In general, style transfer can be divided into two categories: *photorealistic* and *artistic*. Photorealistic style transfer [32,29,56,36] aims at applying reference styles on scenes without hurting details and satisfying contradictory objectives. By contrast, artistic style transfer [14,28,20,35,27,30,4] captures style concepts from reference and

modifies color distributions and texture patterns of content images. However, it requires preparing numerous style images in advance, which limits practicality of style transfer. To tackle this issue, L<sub>D</sub>AST allows following textual descriptions to perform *artistic* style transfer and improves the accessibility of VFX design.

**Language-based Image Editing.** The general task of L<sub>D</sub>AST is language-based image editing (LBIE), which also uses language to edit input images. With rule-based instructions and predefined semantic labels, they [7,25] first carry out LBIE but under limited practicality. Inspired by text-to-image generation [40,58,54], previous works [5,33,26,52,31,9,10,11] perform LBIE by conditional GAN, which modifies the properties of objects in the image. In contrast, L<sub>D</sub>AST aims at preserving the scene structure from the content image and performing stylization guided by the style instruction.

**CLIP-guided Optimization.** Recently, based on the powerful visual-linguistic connection of CLIP [44], CLIP-guided image synthesis [39,34] has shown exciting results. StyleCLIP [37] and NADA [12] tweak the latent code of a pre-trained StyleGAN [23] for image editing. Since heavily relying on a pre-trained generator, both are confined to the training domain, and the results can only present limited stylization. CLIPstyler [24] updates the style transfer network for target style patterns from the CLIP alignment. Though supporting arbitrary content images, CLIPstyler still requires hundreds of iterations and takes lots of time with considerable GPU memory, suffering from the efficiency and practicality overhead. Moreover, our experiments show that CLIP poorly captures detailed style patterns from instructions, which is intractable to perform explicit L<sub>D</sub>AST.

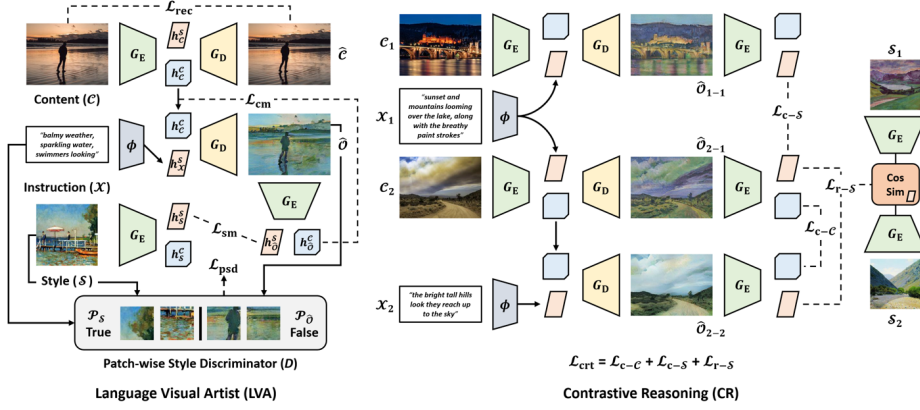
### 3 Language-Driven Artistic Style Transfer

#### 3.1 Overview of CLVA

We introduce language-driven artistic style transfer (L<sub>D</sub>AST) to manipulate the style of a content image  $\mathcal{C}$ , guided by a style instruction  $\mathcal{X}$ , as illustrated in Fig. 1. For training, we have pairs of style images  $\mathcal{S}$  with style instructions  $\mathcal{X}$  to learn the mutual correlation. During testing, only  $\mathcal{X}$  are provided for L<sub>D</sub>AST to carry out artistic style transfer purely relied on language. We present contrastive language visual artist (CLVA) in Fig. 2. Language visual artist (LVA) extracts content structures from  $\mathcal{C}$  and visual patterns from  $\mathcal{X}$  to perform L<sub>D</sub>AST. LVA adopts the patch-wise style discriminator  $D$  to connect extracted visual semantics to patches of paired style image ( $\mathcal{P}_S$  in Fig. 2). Contrastive reasoning (CR) allows comparing contrastive pairs  $\mathcal{C}_1-\mathcal{X}_1$ ,  $\mathcal{C}_2-\mathcal{X}_1$ , and  $\mathcal{C}_2-\mathcal{X}_2$  of content image and style instruction. In this way, it should present consistent content structures from the same content image  $\mathcal{C}_2$  or analogous style patterns from related style images  $\mathcal{S}_1$  and  $\mathcal{S}_2$ , despite using different style instructions.

#### 3.2 Language Visual Artist (LVA)

To tackle L<sub>D</sub>AST, language visual artist (LVA) first adopts visual encoder  $G_E$  to extract the content feature  $h^C$  and the style feature  $h^S$  for an image. Text



**Fig. 2.** Contrastive language visual artist (CLVA). Language Visual Artist (LVA) learns to jointly embed style images  $\mathcal{S}$  and style instructions  $\mathcal{X}$  by the patch-wise style discriminator  $D$  and perform LDATAST for content images  $\mathcal{C}$ . Contrastive Reasoning (CR) compares contrastive pairs to improve the relativeness between transferred results  $\hat{\mathcal{O}}$ .

encoder  $\phi$  also extracts the style instruction feature  $h_{\mathcal{X}}^S$  from an instruction.  $h^C$  is a spatial tensor containing the content structure feature, and  $h^S$  represents the global style pattern.  $\mathcal{S}_{\mathcal{X}}^S$  embeds into the same space of  $h^S$  to reflect the extracted visual semantic. Then, visual decoder  $G_D$  produces transferred results  $\hat{\mathcal{O}}$  from  $h_C^C$  and  $h_{\mathcal{X}}^S$ , which performs style transfer by style instructions:

$$\begin{aligned} h_C^C, h_C^S &= G_E(\mathcal{C}), & h_{\mathcal{X}}^S &= \phi(\mathcal{X}), \\ \hat{\mathcal{O}} &= G_D(h_C^C, h_{\mathcal{X}}^S). \end{aligned} \quad (1)$$

In particular,  $G_D$  applies self-attention [57,35] to fuse  $h^C$  and  $h^S$  over the global spatial dimension. There are two goals to train LVA for LDATAST: (i) preserving *content structures* from content images; (ii) presenting *style patterns* correlated with visual semantics of style instructions.

**Structure Reconstruction.** To preserve content structures, we consider that visual decoder  $G_D$  should be able to reconstruct input content images using extracted content features  $h_C^C$  and style features  $h_C^S$  from visual encoder  $G_E$ :

$$\begin{aligned} \hat{\mathcal{C}} &= G_D(h_C^C, h_C^S), \\ \mathcal{L}_{\text{rec}} &= \|\hat{\mathcal{C}} - \mathcal{C}\|_2, \end{aligned} \quad (2)$$

where the reconstruction loss  $\mathcal{L}_{\text{rec}}$  is computed as the mean L2 difference between reconstructed content images  $\hat{\mathcal{C}}$  and input content images  $\mathcal{C}$ .

**Patch-wise Style Discriminator (D).** Regarding style patterns, results  $\hat{\mathcal{O}}$  guided by style instructions  $\mathcal{X}$  are expected to present analogously to reference style images  $\mathcal{S}$ . To address the connection between linguistic from  $\mathcal{X}$  and visual semantics from  $\mathcal{S}$ , we introduce the patch-wise style discriminator  $D$ . Inspired by texture synthesis [53,15], images with analogous patch patterns should appear



perceptually similar texture patterns.  $D$  tries to recognize the correspondence between an image patch  $\mathcal{P}$  and a style instruction  $\mathcal{X}$ :

$$\begin{aligned}\mathcal{P}_{\hat{\mathcal{O}}}, \mathcal{P}_S &= \text{Crop}(\hat{\mathcal{O}}), \text{Crop}(S), \\ \mathcal{L}_{\text{psd}} &= \log(1 - D(\mathcal{P}_{\hat{\mathcal{O}}}, \mathcal{X})), \\ \mathcal{L}_D &= \log(1 - D(\mathcal{P}_{\hat{\mathcal{O}}}, \mathcal{X})) + \log(D(\mathcal{P}_S, \mathcal{X})),\end{aligned}\tag{3}$$

where **Crop** is to randomly crop an image into patches. The patch-wise style loss  $\mathcal{L}_{\text{psd}}$  aims at generating transferred results that are correlated with  $\mathcal{X}$ . Contrarily, by the discriminator loss  $\mathcal{L}_D$ ,  $D$  learns to distinguish that a patch  $\mathcal{P}$  is from style images ( $\mathcal{P}_S$ ) or transferred results ( $\mathcal{P}_{\hat{\mathcal{O}}}$ ). This adversarial loss [17,41] encourages that transferred results from style instructions are presented similarly with style images, which jointly embeds the extracted visual semantics.

**Content Matching and Style Matching.** To further enhance the alignment with inputs, inspired by cycle consistency [60,51,38,55], we consider the content matching loss  $\mathcal{L}_{\text{cm}}$  and the style matching loss  $\mathcal{L}_{\text{sm}}$  of transferred results  $\hat{\mathcal{O}}$ . We adopt  $G_E$  again to extract content features  $h_{\hat{\mathcal{O}}}^c$  and style features  $h_{\hat{\mathcal{O}}}^s$  for  $\hat{\mathcal{O}}$ , where  $h_{\hat{\mathcal{O}}}^c$  and  $h_{\hat{\mathcal{O}}}^s$  should correlate with  $h_{\mathcal{C}}^c$  from  $\mathcal{C}$  and  $h_{\mathcal{S}}^s$  from  $\mathcal{S}$ :

$$\begin{aligned}(h_{\hat{\mathcal{O}}}^c, h_{\hat{\mathcal{O}}}^s), (h_S^c, h_S^s) &= G_E(\hat{\mathcal{O}}), G_E(S), \\ \mathcal{L}_{\text{cm}}, \mathcal{L}_{\text{sm}} &= \|h_{\hat{\mathcal{O}}}^c - h_{\mathcal{C}}^c\|_2, \|h_{\hat{\mathcal{O}}}^s - h_{\mathcal{S}}^s\|_2.\end{aligned}\tag{4}$$

Therefore, transferred results are required to align with content structures and style patterns from inputs, which meets the goal of LDATAST.

### 3.3 Contrastive Reasoning (CR)

The content image should transfer to various styles while preserving the same structure. Related style instructions can apply analogous style patterns to arbitrary content images. As shown in Fig. 2, contrastive reasoning (CR) compares content structures or style patterns from transferred results of contrastive pairs. The contrastive pair consists of two different content images  $\mathcal{C}_1$  and  $\mathcal{C}_2$  with two reference styles  $\{\mathcal{S}_1, \mathcal{X}_1\}$  and  $\{\mathcal{S}_2, \mathcal{X}_2\}$ . We follow the LVA inference to acquire cross results for pairs of content images and style instructions:

$$\begin{aligned}(h_{\mathcal{C}_1}^c, h_{\mathcal{C}_1}^s), (h_{\mathcal{C}_2}^c, h_{\mathcal{C}_2}^s) &= G_E(\mathcal{C}_1), G_E(\mathcal{C}_2), \\ h_{\mathcal{X}_1}^s, h_{\mathcal{X}_2}^s &= \phi(\mathcal{X}_1), \phi(\mathcal{X}_2), \\ \hat{\mathcal{O}}_{\mathcal{C}_1-\mathcal{X}_1}, \hat{\mathcal{O}}_{\mathcal{C}_1-\mathcal{X}_2} &= G_D(h_{\mathcal{C}_1}^c, h_{\mathcal{X}_1}^s), G_D(h_{\mathcal{C}_1}^c, h_{\mathcal{X}_2}^s), \\ \hat{\mathcal{O}}_{\mathcal{C}_2-\mathcal{X}_1}, \hat{\mathcal{O}}_{\mathcal{C}_2-\mathcal{X}_2} &= G_D(h_{\mathcal{C}_2}^c, h_{\mathcal{X}_1}^s), G_D(h_{\mathcal{C}_2}^c, h_{\mathcal{X}_2}^s).\end{aligned}\tag{5}$$

**Consistent Matching.** Transferred results should present similar content structures ( $\hat{\mathcal{O}}_{\mathcal{C}_2-\mathcal{X}_1}$  and  $\hat{\mathcal{O}}_{\mathcal{C}_2-\mathcal{X}_2}$ ) or analogous style patterns ( $\hat{\mathcal{O}}_{\mathcal{C}_1-\mathcal{X}_1}$  and  $\hat{\mathcal{O}}_{\mathcal{C}_2-\mathcal{X}_1}$ ) if

**Algorithm 1** Training Process of Language Visual Artist (LVA)

---

```

1:  $G_E, G_D$ : Visual Encoder, Visual Decoder
2:  $\phi$ : Text Encoder
3:  $D$ : Patch-wise Style Discriminator
4: while TRAIN_VLA do
5:    $\mathcal{C}, \{\mathcal{S}, \mathcal{X}\} \leftarrow$  Sampled content/style
6:
7:    $h_C^C, h_C^S \leftarrow G_E(\mathcal{C})$     $\hat{\mathcal{C}} \leftarrow G_D(h_C^C, h_C^S)$ 
8:    $\mathcal{L}_{\text{rec}} \leftarrow$  Reconstruction loss ▷ Eq. 2
9:    $h_X^S \leftarrow \phi(\mathcal{X})$     $\hat{\mathcal{O}} \leftarrow G_D(h_C^C, h_X^S)$ 
10:   $\mathcal{P}_S, \mathcal{P}_{\hat{\mathcal{O}}} \leftarrow$  Crop( $\mathcal{S}$ ), Crop( $\hat{\mathcal{O}}$ )
11:   $\mathcal{L}_{\text{psd}} \leftarrow$  Patch-wise style loss ▷ Eq. 3
12:   $(h_{\hat{\mathcal{O}}}^C, h_{\hat{\mathcal{O}}}^S), (h_S^C, h_S^S) \leftarrow G_E(\hat{\mathcal{O}}), G_E(\mathcal{S})$ 
13:   $\mathcal{L}_{\text{cm}} \leftarrow$  Content matching loss ▷ Eq. 4
14:   $\mathcal{L}_{\text{sm}} \leftarrow$  Style matching loss ▷ Eq. 4
15:
16:   $\mathcal{L}_G \leftarrow \mathcal{L}_{\text{rec}} + \mathcal{L}_{\text{psd}} + \mathcal{L}_{\text{cm}} + \mathcal{L}_{\text{sm}}$  ▷ Eq. 3
17:   $\mathcal{L}_D \leftarrow$  Discriminator loss for D ▷ Eq. 3
18:  Update  $G_E, G_D, \phi$  by minimizing  $\mathcal{L}_G$ 
19:  Update  $D$  by maximizing  $\mathcal{L}_D$ 
20: end while

```

---

using the same content image ( $\mathcal{C}_2$ ) or the same style instruction ( $\mathcal{X}_1$ ):

$$\begin{aligned}
h_{\hat{\mathcal{O}}_{\mathcal{C}_i-\mathcal{X}_j}}^C &= G_E(\hat{\mathcal{O}}_{\mathcal{C}_i-\mathcal{X}_j}), \\
\mathcal{L}_{\mathcal{C}-\mathcal{C}} &= \|h_{\hat{\mathcal{O}}_{\mathcal{C}_1-\mathcal{X}_1}}^C - h_{\hat{\mathcal{O}}_{\mathcal{C}_1-\mathcal{X}_2}}^C\|_2 + \|h_{\hat{\mathcal{O}}_{\mathcal{C}_2-\mathcal{X}_1}}^C - h_{\hat{\mathcal{O}}_{\mathcal{C}_2-\mathcal{X}_2}}^C\|_2, \\
\mathcal{L}_{\mathcal{C}-\mathcal{S}} &= \|h_{\hat{\mathcal{O}}_{\mathcal{C}_1-\mathcal{X}_1}}^S - h_{\hat{\mathcal{O}}_{\mathcal{C}_2-\mathcal{X}_1}}^S\|_2 + \|h_{\hat{\mathcal{O}}_{\mathcal{C}_1-\mathcal{X}_2}}^S - h_{\hat{\mathcal{O}}_{\mathcal{C}_2-\mathcal{X}_2}}^S\|_2,
\end{aligned} \tag{6}$$

where *consistent matching* of content structure  $\mathcal{L}_{\mathcal{C}-\mathcal{C}}$  or style pattern  $\mathcal{L}_{\mathcal{C}-\mathcal{S}}$  is aligned by content features or style features, extracted by  $G_E$ .

**Relative Matching.** Apart from consistent matching, distinct style instructions, which imply corresponding visual semantics, should still present relative patterns. For example, we can only discover “*reach up to the sky*” literally from  $\mathcal{X}_2$ . If comparing reference style images  $\mathcal{S}_1$  and  $\mathcal{S}_2$ , we can perceive the sharing of a similar style pattern and link the visual concept of “*bright tall hills*” in  $\mathcal{X}_2$  to “*mountains looming over the lake*” in  $\mathcal{X}_1$ . We define *relative matching*  $\mathcal{L}_{\mathcal{R}-\mathcal{S}}$  with the cosine similarity (CosSim) between reference style images:

$$\begin{aligned}
(h_{\mathcal{S}_i}^C, h_{\mathcal{S}_i}^S) &= G_E(\mathcal{S}_i), \\
r &= \text{CosSim}(h_{\mathcal{S}_1}^S, h_{\mathcal{S}_2}^S), \\
\mathcal{L}_{\mathcal{R}-\mathcal{S}} &= (\|h_{\hat{\mathcal{O}}_{\mathcal{C}_1-\mathcal{X}_1}}^S - h_{\hat{\mathcal{O}}_{\mathcal{C}_1-\mathcal{X}_2}}^S\|_2 + \\
&\quad \|h_{\hat{\mathcal{O}}_{\mathcal{C}_2-\mathcal{X}_1}}^S - h_{\hat{\mathcal{O}}_{\mathcal{C}_2-\mathcal{X}_2}}^S\|_2) \cdot r.
\end{aligned} \tag{7}$$

When style images are related, it has to align style features to certain extent even if paired style instructions are different. Otherwise,  $\mathcal{L}_{\mathcal{R}-\mathcal{S}}$  will be close to 0 and ignore this unrelated style pair. The overall contrastive reasoning loss  $\mathcal{L}_{\text{ctr}}$  considers both consistent matching and relative matching:

$$\mathcal{L}_{\text{ctr}} = \mathcal{L}_{\mathcal{C}-\mathcal{C}} + \mathcal{L}_{\mathcal{C}-\mathcal{S}} + \mathcal{L}_{\mathcal{R}-\mathcal{S}}. \tag{8}$$

### 3.4 Learning of CLVA

For each epoch of CLVA training, we first train with the LVA process and then CR. As algo. 1, we consider reconstruction loss  $\mathcal{L}_{\text{rec}}$  to preserve content structure and patch-wise style loss  $\mathcal{L}_{\text{psd}}$  between style instruction and visual pattern of transferred results. Both content matching loss  $\mathcal{L}_{\text{cm}}$  and style matching loss  $\mathcal{L}_{\text{sm}}$  enhance the matching with the inputs. Simultaneously, we update  $D$  by maximizing discriminator loss  $\mathcal{L}_D$  to distinguish between true patches  $\mathcal{P}_S$  or false patches  $\mathcal{P}_{\hat{O}}$ , concerning style instructions. During CR, contrastive pairs of content images and style instructions are randomly sampled, and the transferred results are across produced. We further update by minimizing contrastive reasoning loss  $\mathcal{L}_{\text{ctr}}$  to allow considering content consistency and mutual style relativeness. The overall optimization of CLVA is summarized as:

$$\begin{aligned} \mathcal{L}_G &= \mathcal{L}_{\text{rec}} + \mathcal{L}_{\text{psd}} + \mathcal{L}_{\text{cm}} + \mathcal{L}_{\text{sm}}, \\ \min_{G, \phi} \max_D \mathcal{L}_G + \mathcal{L}_D + \mathcal{L}_{\text{ctr}}. \end{aligned} \quad (9)$$

## 4 Experiments

### 4.1 Experimental Setup

**Dataset.** To evaluate our CLVA, we consider DTD<sup>2</sup> [8] and ArtEmis [1] as reference style instructions. DTD<sup>2</sup> contains 5K texture images with its natural descriptions for visual attributes such as colors and texture patterns. ArtEmis provides 80K artworks from WikiArt<sup>1</sup> with annotations of visual contents and emotional effects as human style feelings. We also collect 15K wallpapers from WallpapersCraft<sup>2</sup>, which presents diverse scenes as content images. Each content image is resized to 256x192 in our experiment. We randomly sample 100 unseen content images and 100 testing reference styles to evaluate the generalizability of LDATA. Note that both style images and style instructions appear for training, but only style instructions are accessible during testing.

**Evaluation Metrics.** To support large-scale evaluation, we treat transferred results directly from style images as semi-groundtruth (Semi-GT) [2,3,42] by the SOTA style transfer AdaAttn [30]. We apply the following metrics:

- **SSIM** [48] compares images in the luminance, contrast, and structure aspects. A higher SSIM has a higher structural similarity;
- **Percept** [22] computes from the gram matrix of visual features. A lower Percept loss shows that two images share a similar style pattern;
- **FAD** [18] is computed by the mean L2 distance of the activations from the InceptionV3 [46] feature. As a distance metric, a lower FAD represents that LDATA results and Semi-GT are more relevant.

<sup>1</sup> WikiArt: <https://www.wikiart.org>

<sup>2</sup> WallpapersCraft: <https://wallpaperscraft.com/>

Note that we consider SSIM and FAD to compare with Semi-GT and calculate Percept loss directly with reference style images. Apart from visual similarity, we consider the correlation between style instructions and LDAST results:

- **VLS** [49] calculates the cosine similarity between each other from CLIP [44]. Since each metric has different deficiencies, we also conduct a comprehensive human evaluation from aspects of content, instruction, and style matching. We randomly sample 75 LDAST results and adopt MTurk<sup>3</sup> to rank over all methods. We also hire 3 MTurkers for each task to avoid the potential ranking bias.

**Baselines.** We conduct baselines for LDAST from various aspects:

- **Style Transfer:** We consider previous artistic style transfer methods NST [14], WCT [28], AdaIn [20], SANet [35], and LST [27] that support arbitrary content images. We use the same style (instruction and image) encoding from our CLVA as style features and follow their own training process to perform LDAST upon them. Due to the space issue, we only show the comparison with more recent SANet and LST. Please refer to Appendix for the complete results.
- **Language-based Image Editing:** We adopt ManiGAN [26] with affine combination module (ACM) as the general language-based editing baseline, where it modifies the content image by the style instruction. We treat normal style transferred results as groundtruth for ManiGAN to learn from.
- **CLIP-based Optimization:** StyleCLIP [37], NADA [12], and CLIPstyler [24] manipulate the content image based on the CLIP alignment of the guided instruction. Since StyleCLIP and NADA are restricted by the pre-trained generator, we compare them with the training domains of car and church. Differently, CLIPstyler can carry out arbitrary content images for LDAST.

## 4.2 Quantitative Results

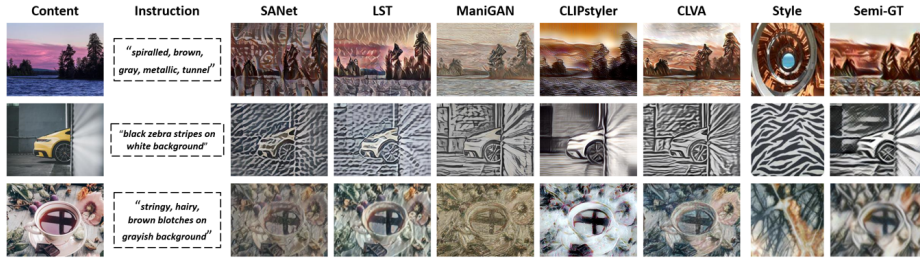
**Instruction with Visual Attributes.** Table 1 illustrates the comparison of LDAST with baselines on DTD<sup>2</sup>. As regards automatic metrics, CLVA preserves content structures (highest 36.65 SSIM) and stylizes with related visual attributes to style images (lowest 0.2033 Percept loss). Furthermore, CLVA brings out the highest overall similarity as Semi-GT (lowest 0.1493 FAD). Since CLIPstyler directly optimizes by CLIP [44], it makes the highest VLS. Through the patch-wise discriminator, our CLVA can still produce style patterns correlated to given instructions (competitive 24.00 VLS) even without the pre-trained CLIP.

The human evaluation investigates the matching between transferred results with content images (Content), style instructions (Instruction), style images (Style), and Semi-GT (Semi-GT). In particular, content and instruction matching are the two most crucial, which concern the goal of LDAST: *content structure preservation* and *style pattern presentation*; style image and semi-gt matching are provided for different comparing targets from a human aspect. The results are calculated by the mean ranking score (from 1 to 5, the higher is better) of each method. In general, MTurkers indicate that our CLVA has an apparent advantage in preserving content structures (highest 3.852 Content) and presenting aligned

<sup>3</sup> Amazon Mechanical Turk: <https://www.mturk.com>

| Method          | Automatic Metrics |               |               |              | Human Evaluation |              |              |              |
|-----------------|-------------------|---------------|---------------|--------------|------------------|--------------|--------------|--------------|
|                 | SSIM↑             | Percept↓      | FAD↓          | VLS↑         | Content↑         | Instruction↑ | Style↑       | Semi-GT↑     |
| SANet [35]      | 35.50             | 0.2129        | 0.1627        | 23.57        | 2.701            | 2.477        | 2.738        | 2.630        |
| LST [27]        | <u>34.84</u>      | 0.2129        | <u>0.1533</u> | 23.16        | 2.743            | 2.831        | 2.651        | 2.528        |
| ManiGAN [26]    | 32.70             | 0.2401        | 0.1663        | 23.25        | 2.757            | 2.562        | 2.937        | 2.922        |
| CLIPstyler [24] | 25.24             | 0.2598        | 0.1818        | <b>24.62</b> | <u>2.948</u>     | <u>3.388</u> | <u>3.073</u> | <u>3.265</u> |
| CLVA            | <b>36.65</b>      | <b>0.2033</b> | <b>0.1493</b> | <u>24.00</u> | <b>3.852</b>     | <b>3.742</b> | <b>3.603</b> | <b>3.655</b> |

**Table 1.** Testing results of LDASt using visual attribute instructions on DTD<sup>2</sup>.



**Fig. 3.** Visualized comparison using visual attribute instructions on DTD<sup>2</sup>.

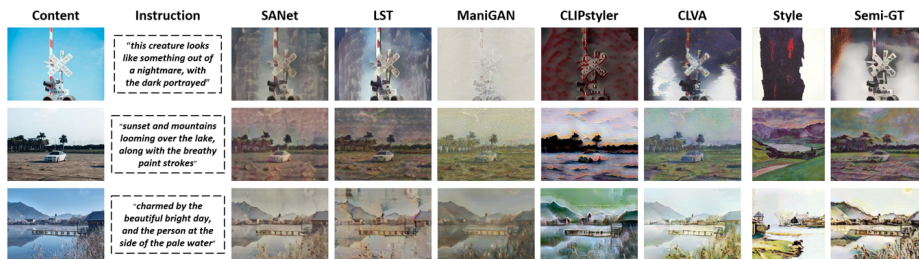
style patterns (highest 3.742 Instruction). Though with the aid of CLIP, CLIPstyler is still behind CLVA (-0.4 Instruction), with an even higher gap in style image matching (-0.5 Style). Contributed by contrastive reasoning that compares the mutual relativeness between pairs of contents and instructions, CLVA can stylize with the captured visual attributes. We adopt Pearson correlation and investigate the coefficients between automatic metrics and human evaluation as 77.2 (FAD→Instruction), 84.5 (FAD→Semi-GT), 81.3 (VLS→Instruction), and 77.8 (VLS→Semi-GT). This high correlation indicates that our metric design is adequate for evaluating large-scale LDASt experiments. The even higher 88.2 correlation (Instruction→Semi-GT) between instruction and Semi-GT matching in human evaluation further supports the usage of Semi-GT.

From the aspect of visualized comparison in Fig. 3, previous SANet and LST only produce repetitive and disorder textures in their transferred results. ManiGAN modifies the style directly over pixels, suffering from blurring objects; this deficiency can also be found in Table 1 (lower SSIM and lower Content matching). CLIPstyler is sometimes misguided by CLIP, making irrelevant patterns, such as the bright white background in the third case. Contrary to baselines, CLVA extracts a more detailed style from different kinds of guidance (“*brown metallic*” in the first row and “*stringy hairy*” in the third case), leading to superior LDASt results that correspond to style instructions.

**Instruction with Emotional Effects.** Unlike visual attributes, emotional effect instructions are more challenging as connecting to visual semantics of described objects or style patterns from human feelings. For example, “*yellowish and green*” from “*sunset and mountains*” or “*scaring charcoal grey*” from “*nightmare*”. We consider this human style feeling on ArtEmis [1], where the model has to express the latent visual concepts of emotional effect instructions. CLVA performs with more balance (both second-highest SSIM and second-lowest Per-

| Method          | Automatic Metrics |               |               |              | Human Evaluation |              |              |              |
|-----------------|-------------------|---------------|---------------|--------------|------------------|--------------|--------------|--------------|
|                 | SSIM↑             | Percept↓      | FAD↓          | VLS↑         | Content↑         | Instruction↑ | Style↑       | Semi-GT↑     |
| SANet [35]      | 38.36             | <b>0.0352</b> | 0.1548        | 19.30        | 3.170            | 2.978        | 2.980        | 2.890        |
| LST [27]        | <b>42.13</b>      | 0.0386        | 0.1595        | 19.92        | 2.967            | 2.714        | 2.614        | 2.757        |
| ManiGAN [26]    | 38.46             | 0.0500        | 0.1554        | 19.69        | 2.729            | 2.583        | 2.879        | <b>3.192</b> |
| CLIPstyler [24] | 24.17             | 0.0659        | 0.1759        | <b>21.04</b> | 2.777            | <b>3.140</b> | <b>2.998</b> | 2.952        |
| CLVA            | <b>40.32</b>      | <b>0.0357</b> | <b>0.1418</b> | <b>20.11</b> | <b>3.357</b>     | <b>3.586</b> | <b>3.530</b> | <b>3.208</b> |

**Table 2.** Testing results of LDAST using emotional effect instructions on ArtEmis.



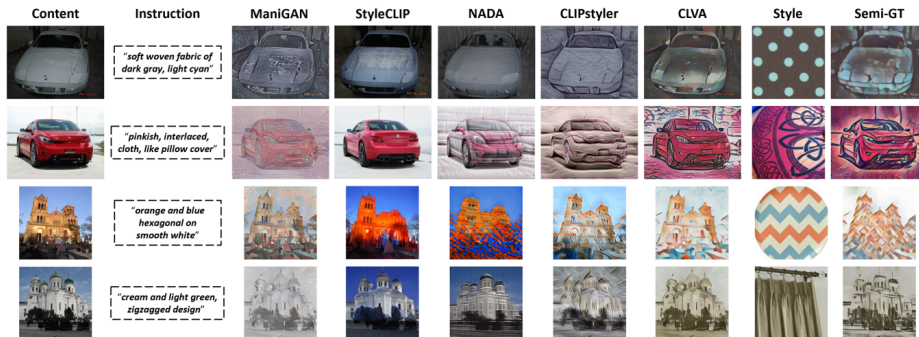
**Fig. 4.** Visualized comparison using emotional effect instructions on ArtEmis.

cept) from Table 2, especially the lowest 0.1418 FAD, making the most similar transferred results to Semi-GT. Though CLIPstyler [24] achieves higher VLS by optimizing over CLIP, from human aspects, CLVA can preserve more concrete contents and present more correlated style patterns (higher 3.357 content and 3.586 instruction matching). The visualized comparison in Fig. 4 illustrates that previous SANet [35] and LST [27] contain unsmooth and fragmentary patterns with blurring contents. Without a style transformation process, ManiGAN [26] modifies with only monotonous colors. CLIPstyler is failed to capture human style feelings well, suffering from weird and unpleasant results. Different from them, our CLVA learns the visual semantic during contrastive reasoning by comparing mutual relativeness between literal instructions and style images, leading to a more colorful and corresponding stylization as human emotion. More surprisingly, despite not instructed literally, CLVA perceives “*side of the water*” and reveals the latent yet correlated “*grassland*” precisely in the third row.

**Specific Content Domain.** To compare with StyleCLIP [37] and NADA [12] that are restricted by the pre-trained generator, we evaluate LDAST on the specific content domain. We consider the same domain images in StyleGAN2 [23] and visual attribute instructions on DTD<sup>2</sup>. Table 3 indicates the numerical comparison on Car and Church. Our CLVA still produces superior results and is the most admirable by human. Since StyleCLIP and NADA rely on StyleGAN, they can only preserve content (highest 3.459 Content by StyleCLIP) but with limited stylization (lower Instruction and Style). Similar observations can be found in Fig. 5, where StyleCLIP shows almost no modification for the second car. They can neither deal with the background; NADA even destroys the scene in the third row. In contrast to CLIPstyler [24] that only contains abstractive and obscure styles, CLVA presents the detailed “*read interplaced cloth*” behind the car and the color “*cream*” precisely on the surface of the church.

| Method          | Automatic Metrics |               |               |              | Human Evaluation |              |              |              |
|-----------------|-------------------|---------------|---------------|--------------|------------------|--------------|--------------|--------------|
|                 | SSIM↑             | Percept↓      | FAD↓          | VLS↑         | Content↑         | Instruction↑ | Style↑       | Semi-GT↑     |
| ManiGAN [26]    | 26.45             | 0.2329        | 0.1672        | 23.44        | 2.861            | 2.894        | 2.978        | 2.893        |
| StyleCLIP [37]  | <u>28.03</u>      | 0.2609        | 0.1812        | 21.55        | <b>3.459</b>     | 2.845        | 2.930        | 2.829        |
| NADA [12]       | 16.98             | 0.2733        | 0.1876        | 23.38        | 2.542            | 2.798        | 2.846        | 2.932        |
| CLIPstyler [24] | 18.43             | 0.2493        | 0.1826        | <b>24.16</b> | 2.986            | <b>3.067</b> | <b>3.003</b> | <b>3.032</b> |
| CLVA            | <b>30.98</b>      | <b>0.1957</b> | <b>0.1544</b> | <u>23.68</u> | <u>3.153</u>     | <b>3.465</b> | <b>3.344</b> | <b>3.315</b> |

**Table 3.** Testing results of LDATA on specific content domain (Car and Church).



**Fig. 5.** Visualized comparison on specific content domain (Car and Church).

### 4.3 Ablation Study

We conduct an ablation study of each component effect on DTD<sup>2</sup> in Table 4. At row (a), with the reconstruction  $\mathcal{L}_{\text{rec}}$  and the patch-wise style  $\mathcal{L}_{\text{psd}}$ , CLVA achieves feasible LDATA results by concrete structures and extracted style semantics. Row (b)-(d) shows the strength of content matching  $\mathcal{L}_{\text{cm}}$  and style matching  $\mathcal{L}_{\text{sm}}$ . In particular, content matching helps the structure similarity to content images (higher 36.05 SSIM). Style matching aims at analogous visual patterns to style images, which leads to better stylization quality (lower 0.2049 Percept and higher 23.69 VLS). If considering altogether, it can benefit and strike a balance between both. Finally, contrastive reasoning  $\mathcal{L}_{\text{ctr}}$  further enables CLVA to consider contrastive pairs, making a comprehensive improvement at row (e).

**Why CLVA is better than CLIP-based?** Despite no CLIP optimized, CLVA demonstrates superior results on LDATA with all aspects of automatic metrics and human evaluation. To investigate it, we conduct instruction-to-style retrieval based on the similarity between features of style instructions and style images. Table 5 shows that our learned CLVA performs higher Recall@k on both DTD<sup>2</sup> and ArtEmis, leading to a better instruction-style alignment than the used CLIP. The visualization in Fig. 6 also indicates the flaw of CLIP on detailed style patterns. For example, in the first row, CLIP only presents either “*bright color*” or “*town*” in the retrieval results. In contrast, CLVA can capture both and present more related LDATA to “*happy place to live*”. From Table 6, even CLIP has been fine-tuned ahead; our CLVA still produces preferable LDATA results from all human aspects of content, instruction, and style matching. This observation supports that contrastive reasoning, which considers contrastive pairs of content images and style instructions, is required to benefit from mutual relativeness.



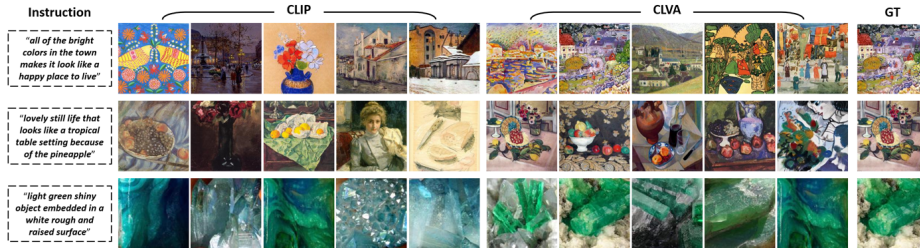
|     | Ablation Settings                                     |                           |                           |                            | Automatic Metrics |                      |                  |                |
|-----|---|---------------------------|---------------------------|----------------------------|-------------------|----------------------|------------------|----------------|
|     | $\mathcal{L}_{\text{rec}} + \mathcal{L}_{\text{psd}}$ | $\mathcal{L}_{\text{cm}}$ | $\mathcal{L}_{\text{sm}}$ | $\mathcal{L}_{\text{ctr}}$ | SSIM $\uparrow$   | Percept $\downarrow$ | FAD $\downarrow$ | VLS $\uparrow$ |
| (a) | ✓   | ✗                         | ✗                         | ✗                          | 34.73             | 0.2290               | 0.1568           | 23.29          |
| (b) | ✓   | ✓                         | ✗                         | ✗                          | <u>36.05</u>      | 0.2304               | 0.1512           | 23.27          |
| (c) | ✓   | ✗                         | ✓                         | ✗                          | 35.73             | <u>0.2049</u>        | <u>0.1508</u>    | 23.69          |
| (d) | ✓   | ✓                         | ✓                         | ✗                          | 35.86             | 0.2100               | 0.1499           | 23.54          |
| (e) | ✓   | ✓                         | ✓                         | ✓                          | <b>36.65</b>      | <b>0.2033</b>        | <b>0.1493</b>    | <b>24.00</b>   |

**Table 4.** Ablation study of CLVA using visual attribute instructions on DTD<sup>2</sup>.

| Method    | DTD <sup>2</sup> |             | ArtEmis     |             | Method           | Human Evaluation   |                        |                  |                    |
|-----------|------------------|-------------|-------------|-------------|------------------|--------------------|------------------------|------------------|--------------------|
|           | R@1              | R@5         | R@1         | R@5         |                  | Content $\uparrow$ | Instruction $\uparrow$ | Style $\uparrow$ | Semi-GT $\uparrow$ |
| CLIP [44] | 13.9             | 30.7        | 9.8         | 20.7        | CLIPstyler (ft.) | 1.208              | 1.347                  | 1.292            | 1.333              |
| CLVA      | <b>19.3</b>      | <b>45.1</b> | <b>13.9</b> | <b>30.7</b> | CLVA             | <b>1.792</b>       | <b>1.653</b>           | <b>1.708</b>     | <b>1.667</b>       |

**Table 5.** Instruction-to-style retrieval on DTD<sup>2</sup> and ArtEmis.

**Table 6.** Human comparison between CLVA and CLIPstyle with fine-tuned CLIP on DTD<sup>2</sup>.



**Fig. 6.** Visualization examples of instruction-to-style retrieval by CLIP and CLVA.

Apart from transfer quality, CLVA also holds a higher efficiency than CLIP-based methods. Table 7 illustrates the time and GPU cost on a single TITAN X (12GB) with content image size 256x192. All CLIP-based methods take more than 30 seconds for only one pair of content images and style instructions. Instead of numerous iterations to align with CLIP, we extract style semantics and carry out L<sub>DA</sub>ST in one shot, taking merely 0.03 seconds for one input. Without updating the model during inference, CLVA supports parallelization and can accomplish 50 pairs in half a second. Besides, as a lightweight style transfer network, CLVA requires the least GPU memory for L<sub>DA</sub>ST. In summary, our CLVA surpasses those CLIP-based methods on both quality and efficiency because of the detailed style deficiency and the required optimizing iteration from CLIP.

**Qualitative Results.** As shown in Fig. 7, we investigate the linear interpolation of extracted style patterns by CLVA. Considering style features  $h_{\mathcal{X}_1}^S$  and  $h_{\mathcal{X}_2}^S$  of instructions  $\mathcal{X}_1$  and  $\mathcal{X}_2$ , the interpolated  $h_p^S$  should be:

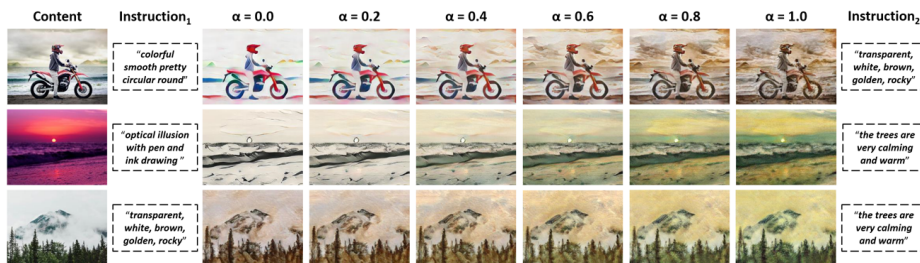
$$h_p^S = (1 - \alpha)h_{\mathcal{X}_1}^S + \alpha h_{\mathcal{X}_2}^S, \quad (10)$$

where  $\alpha$  is the style ratio between the two. Fig. 7 presents a smooth transformation from one style instruction to another. By training on DTD<sup>2</sup> and ArtEmis altogether, CLVA even performs interpolated stylization by both visual attribute and emotional effect instructions in the third row. Fig. 8 illustrates diverse L<sub>DA</sub>ST



| Method          | Time (sec)   |              |              | GPU (MB)    |             |             |
|-----------------|--------------|--------------|--------------|-------------|-------------|-------------|
|                 | BS=1         | 32           | 50           | BS=1        | 32          | 50          |
| ManiGAN [26]    | 0.079        | 0.533        | 1.148        | 3312        | 6572        | 8129        |
| StyleCLIP [37]  | 32.38        | *            | *            | 4149        | *           | *           |
| NADA [12]       | 63.49        | *            | *            | 6413        | *           | *           |
| CLIPstyler [24] | 99.98        | *            | *            | 5429        | *           | *           |
| CLVA            | <b>0.029</b> | <b>0.246</b> | <b>0.405</b> | <b>1525</b> | <b>3207</b> | <b>4441</b> |

**Table 7.** Time and GPU cost when performing LDAST on TITAN X with content image size 256x192. \* means this method can only run one input at a time.



**Fig. 7.** Style interpolation results of LDAST over instructions.

results by our CLVA. Since CLVA supports arbitrary content images, we can also modify the style detail for high-resolution inputs in Fig. 9.

## 5 Conclusion

We introduce language-driven artistic style transfer (LDAST) to do stylization for a content image by a style instruction. We propose contrastive language visual artist (CLVA) that adopts the patch-wise style discriminator and contrastive reasoning to jointly learn between style images and style instructions. We demonstrate that CLVA can express various style patterns of visual attributes as well as emotional effects and perform LDAST efficiently. CLVA also outperforms baselines on both automatic metrics and human evaluation. We believe that LDAST can make visual applications like image/video effect more controllable for humans.

**Acknowledgments.** Research was sponsored by the U.S. Army Research Office and was accomplished under Contract Number W911NF-19-D-0001 for the Institute for Collaborative Biotechnologies. The views and conclusions contained in this document are those of the authors and should not be interpreted as representing the official policies, either expressed or implied, of the U.S. Government. The U.S. Government is authorized to reproduce and distribute reprints for Government purposes notwithstanding any copyright notation herein.

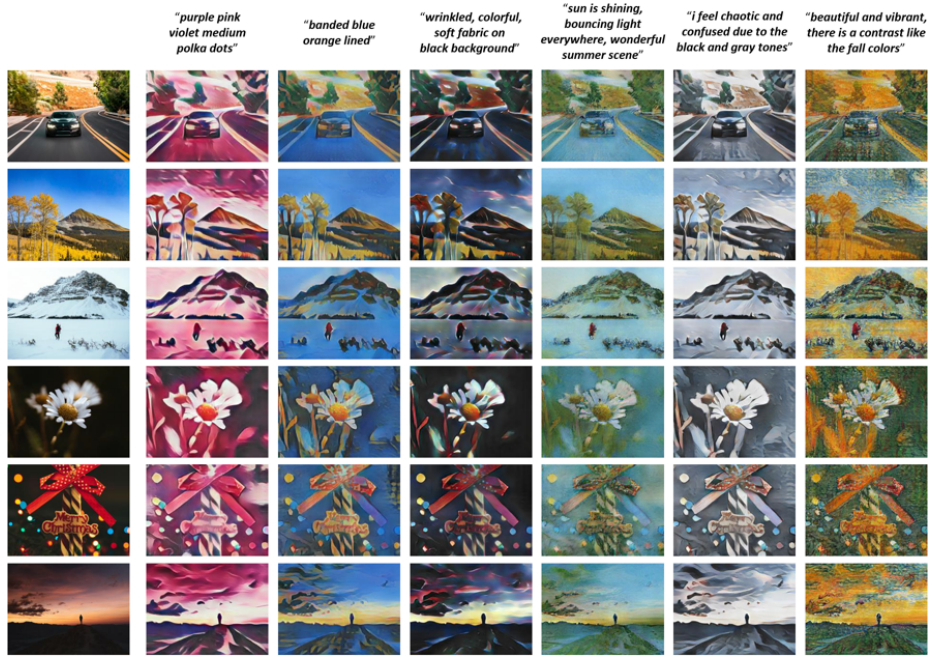


Fig. 8. CLVA results on diverse pairs of content images and style instructions.



Fig. 9. High-resolution (1920x1080) LDAST results by CLVA with upper right: *“the lonely world makes me feel scared and nostalgic how sky and sea merge together”*; lower left: *“the snow and lights in the shop windows looks like a winter scene”*; lower right: *“ink painting, black dotted line, whiteboard”*.

## References

1. Achlioptas, P., Ovsjanikov, M., Haydarov, K., Elhoseiny, M., Guibas, L.: ArtEmis: Affective Language for Visual Art. In: CVPR (2021)
2. Al-Sarraf, A., Shin, B.S., Xu, Z., Klette, R.: Ground Truth and Performance Evaluation of Lane Border Detection. In: ICCVG (2014)
3. Borkar, A., Hayes, M., Smith, M.T.: An Efficient Method to Generate Ground Truth for Evaluating Lane Detection Systems. In: ICASSP (2010)
4. Chen, H., Zhao, L., Wang, Z., Zhang, H., Zuo, Z., Li, A., Xing, W., Lu, D.: DualAST: Dual Style-Learning Networks for Artistic Style Transfer. In: CVPR (2021)
5. Chen, J., Shen, Y., Gao, J., Liu, J., Liu, X.: Language-Based Image Editing with Recurrent Attentive Models. In: CVPR (2018)
6. Chen, Y.L., Hsu, C.T.: Towards Deep Style Transfer: A Content-Aware Perspective. In: BMVC (2016)
7. Cheng, M.M., Zheng, S., Lin, W.Y., Warrell, J., Vineet, V., Sturges, P., Crook, N., Mitra, N., Torr, P.: ImageSpirit: Verbal Guided Image Parsing. In: ACM Transactions on Graphics (2013)
8. Cimpoi, M., Maji, S., Kokkinos, I., Mohamed, S., Vedaldi, A.: Describing Textures in the Wild. In: CVPR (2014)
9. El-Nouby, A., Sharma, S., Schulz, H., Hjelm, D., Asri, L.E., Kahou, S.E., Bengio, Y., W.Taylor, G.: Tell, Draw, and Repeat: Generating and Modifying Images Based on Continual Linguistic Instruction. In: ICCV (2019)
10. Fu, T.J., Wang, X.E., Grafton, S., Eckstein, M., Wang, W.Y.: SSCR: Iterative Language-Based Image Editing via Self-Supervised Counterfactual Reasoning. In: EMNLP (2020)
11. Fu, T.J., Wang, X.E., Grafton, S., Eckstein, M., Wang, W.Y.: M3L: Language-based Video Editing via Multi-Modal Multi-Level Transformer. In: CVPR (2022)
12. Gal, R., Patashnik, O., Maron, H., Chechik, G., Cohen-Or, D.: StyleGAN-NADA: CLIP-Guided Domain Adaptation of Image Generators. In: arXiv:2108.00946 (2021)
13. Gao, C., Gu, D., Zhang, F., Yu, Y.: ReCoNet: Real-time Coherent Video Style Transfer Network. In: arXiv:1807.01197 (2018)
14. Gatys, L.A., Ecker, A.S., Bethge, M.: A Neural Algorithm of Artistic Style. In: arXiv:1508.06576 (2015)
15. Gatys, L.A., Ecker, A.S., Bethge, M.: Texture Synthesis Using Convolutional Neural Networks. In: NeurIPS (2015)
16. Gatys, L.A., Ecker, A.S., Bethge, M., Hertzmann, A., Shechtman, E.: Controlling Perceptual Factors in Neural Style Transfer. In: CVPR (2017)
17. Goodfellow, I.J., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., Courville, A., Bengio, Y.: Generative Adversarial Networks. In: NeurIPS (2014)
18. Heusel, M., Ramsauer, H., Unterthiner, T., Nessler, B., Hochreiter, S.: GANs Trained by a Two Time-Scale Update Rule Converge to a Local Nash Equilibrium. In: NeurIPS (2017)
19. Huang, H., Wang, H., Luo, W., Ma, L., Jiang, W., Zhu, X., Li, Z., Liu, W.: Real-Time Neural Style Transfer for Videos. In: CVPR (2017)
20. Huang, X., Belongie, S.: Arbitrary Style Transfer in Real-time with Adaptive Instance Normalization. In: ICCV (2017)
21. Jing, Y., Yang, Y., Feng, Z., Ye, J., Yu, Y., Song, M.: Neural Style Transfer: A Review. In: arXiv:1705.04058 (2017)

22. Johnson, J., Alahi, A., Fei-Fei, L.: Perceptual Losses for Real-Time Style Transfer and Super-Resolution. In: ECCV (2016)
23. Karras, T., Laine, S., Aittala, M., Hellsten, J., Lehtinen, J., Aila, T.: Analyzing and Improving the Image Quality of StyleGAN. In: CVPR (2020)
24. Kwon, G., Ye, J.C.: CLIPstyler: Image Style Transfer with a Single Text Condition. In: CVPR (2022)
25. Laput, G., Dontcheva, M., Wilensky, G., Chang, W., Agarwala, A., Linder, J., Adar, E.: PixelTone: A Multimodal Interface for Image Editing. In: CHI (2013)
26. Li, B., Qi, X., Lukasiewicz, T., Torr, P.H.S.: ManiGAN: Text-Guided Image Manipulation. In: CVPR (2020)
27. Li, X., Liu, S., Kautz, J., Yang, M.H.: Learning Linear Transformations for Fast Arbitrary Style Transfer. In: CVPR (2019)
28. Li, Y., Fang, C., Yang, J., Wang, Z., Lu, X., Yang, M.H.: Universal Style Transfer via Feature Transforms. In: NeurIPS (2017)
29. Li, Y., Liu, M.Y., Li, X., Yang, M.H., Kautz, J.: A Closed-form Solution to Photorealistic Image Stylization. In: ECCV (2018)
30. Liu, S., Lin, T., He, D., Li, F., Wang, M., Li, X., Sun, Z., Li, Q., Ding, E.: AdaAttN: Revisit Attention Mechanism in Arbitrary Neural Style Transfer. In: ICCV (2021)
31. Liu, X., Lin, Z., Zhang, J., Zhao, H., Tran, Q., Wang, X., Li, H.: Open-Edit: Open-Domain Image Manipulation with Open-Vocabulary Instructions. In: ECCV (2020)
32. Luan, F., Paris, S., Shechtman, E., Bala, K.: Deep Photo Style Transfer. In: CVPR (2017)
33. Nam, S., Kim, Y., Kim, S.J.: Text-Adaptive Generative Adversarial Networks: Manipulating Images with Natural Language. In: NeurIPS (2018)
34. Nichol, A., Dhariwal, P., Ramesh, A., Shyam, P., Mishkin, P., McGrew, B., Sutskever, I., Chen, M.: GLIDE: Towards Photorealistic Image Generation and Editing with Text-Guided Diffusion Models. In: arXiv:2112.10741 (2021)
35. Park, D.Y., Lee, K.H.: Arbitrary Style Transfer with Style-Attentional Networks. In: CVPR (2019)
36. Park, T., Zhu, J.Y., Wang, O., Lu, J., Shechtman, E., Efros, A.A., Zhang, R.: Swapping Autoencoder for Deep Image Manipulation. In: NeurIPS (2020)
37. Patashnik, O., Wu, Z., Shechtman, E., Cohen-Or, D., Lischinski, D.: StyleCLIP: Text-Driven Manipulation of StyleGAN Imagery. In: ICCV (2021)
38. Qiao, T., Zhang, J., Xu, D., Tao, D.: MirrorGAN: Learning Text-to-image Generation by Redescription. In: CVPR (2019)
39. Ramesh, A., Pavlov, M., Goh, G., Gray, S., Voss, C., Radford, A., Chen, M., Sutskever, I.: Zero-Shot Text-to-Image Generation. In: arXiv:2102.12092 (2021)
40. Reed, S., Akata, Z., Yan, X., Logeswaran, L., Schiele, B., Lee, H.: Generative Adversarial Text to Image Synthesis. In: ICML (2016)
41. Salehi, P., Chalechale, A., Taghizadeh, M.: Generative Adversarial Networks (GANs): An Overview of Theoretical Model, Evaluation Metrics, and Recent Developments. In: arXiv:2005.13178 (2020)
42. Salvo, R.D.: Large Scale Ground Truth Generation for Performance Evaluation of Computer Vision Methods. In: VIGTA (2013)
43. Sanakoyeu, A., Kotovenko, D., Lang, S., Ommer, B.: A Style-Aware Content Loss for Real-time HD Style Transfer. In: ECCV (2018)
44. Shi, L., Shuang, K., Geng, S., Su, P., Jiang, Z., Gao, P., Fu, Z., de Melo, G., Su, S.: Contrastive Visual-Linguistic Pretraining. In: arXiv:2007.13135 (2020)
45. Somavarapu, N., Ma, C.Y., Kira, Z.: Frustratingly Simple Domain Generalization via Image Stylization. In: arXiv:2006.11207 (2020)

46. Szegedy, C., Vanhoucke, V., Ioffe, S., Shlens, J., Wojna, Z.: Rethinking the Inception Architecture for Computer Vision. In: CVPR (2016)
47. Wang, P., Li, Y., Vasconcelos, N.: Rethinking and Improving the Robustness of Image Style Transfer. In: CVPR (2021)
48. Wang, Z., Bovik, A.C., Sheikh, H.R., Simoncel, E.P.: Image Quality Assessment: From Error Visibility to Structural Similarity. In: TIP (2004)
49. Wu, C., Huang, L., Zhang, Q., Li, B., Ji, L., Yang, F., Sapiro, G., Duan, N.: GODIVA: Generating Open-Domain Videos from Natural Descriptions. In: arXiv:2104.14806 (2021)
50. Wu, C., Timm, M., Maji, S.: Describing Textures using Natural Language. In: ECCV (2020)
51. Wu, L., Wang, Y., Shao, L.: Cycle-Consistent Deep Generative Hashing for Cross-Modal Retrieval. In: TIP (2018)
52. Xia, W., Yang, Y., Xue, J.H., Wu, B.: TediGAN: Text-Guided Diverse Face Image Generation and Manipulation. In: CVPR (2021)
53. Xian, W., Sangkloy, P., Agrawal, V., Raj, A., Lu, J., Fang, C., Yu, F., Hays, J.: TextureGAN: Controlling Deep Image Synthesis with Texture Patches. In: CVPR (2018)
54. Xu, T., Zhang, P., Huang, Q., Zhang, H., Gan, Z., Huang, X., Hes, X.: AttnGAN: FineGrained Text to Image Generation with Attentional Generative Adversarial Networks. In: CVPR (2018)
55. Yi, Z., Zhang, H., Tan, P., Gong, M.: DualGAN: Unsupervised Dual Learning for Image-to-Image Translation. In: ICCV (2017)
56. Yoo, J., Uh, Y., Chun, S., Kang, B., Ha, J.W.: Photorealistic Style Transfer via Wavelet Transforms. In: ICCV (2019)
57. Zhang, H., Goodfellow, I., Metaxas, D., Odena, A.: Self-Attention Generative Adversarial Networks. In: PMLR (2019)
58. Zhang, H., Xu, T., Li, H., Zhang, S., Wang, X., Huang, X., Metaxas, D.: StackGAN: Text to Photo-realistic Image Synthesis with Stacked Generative Adversarial Networks. In: ICCV (2017)
59. Zhang, Z., Wang, Z., Lin, Z., Qi, H.: Image Super-Resolution by Neural Texture Transfer. In: CVPR (2019)
60. Zhu, J.Y., Park, T., Isola, P., Efros, A.A.: Unpaired Image-to-Image Translation using Cycle-Consistent Adversarial Networks. In: ICCV (2017)