# Crossing You in Style:
# Cross-modal Style Transfer from Music to Visual Arts

Cheng-Che Lee*
nctusunnerli.cs06g@nctu.edu.tw
National Chiao Tung University
Hsinchu, Taiwan

Wan-Yi Lin
softcat477@gmail.com
National Tsing Hua University
Hsinchu, Taiwan

Yen-Ting Shih
steven88sky@gapp.nthu.edu.tw
National Tsing Hua University
Hsinchu, Taiwan

Pei-Yi Patricia Kuo
pykuo@iss.nthu.edu.tw
National Tsing Hua University
Hsinchu, Taiwan

Li Su
lisu@iis.sinica.edu.tw
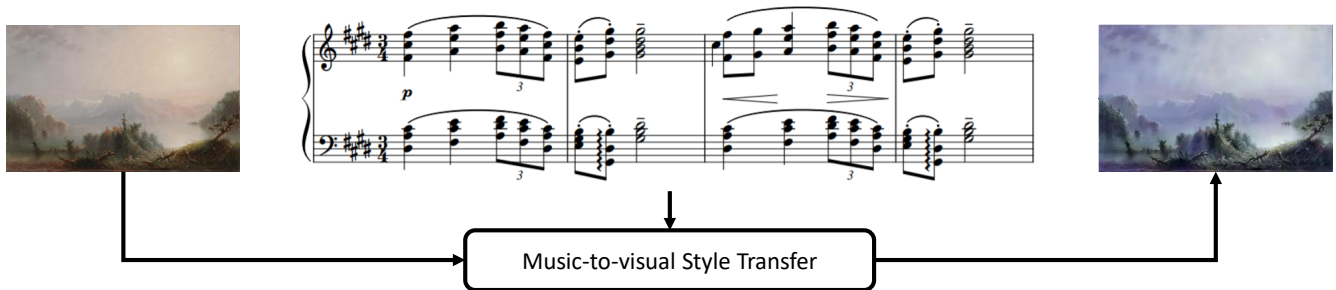Academia Sinica
Taipei, Taiwan

**Figure 1: Let music change the visual style of an image. For example, a spectrogram of Claude Debussy's music *Sarabande* in *Pour le piano, L. 95* (1901) transfers Alfred Jacob Miller's painting *The Lake Her Lone Bosom Expands to the Sky* (1850) into an Impressionism-like color scheme through a neural network linking the semantic space shared by music and image.**

## ABSTRACT

Music-to-visual style transfer is a challenging yet important cross-modal learning problem in the practice of creativity. Its major difference from the traditional image style transfer problem is that the style information is provided by music rather than images. Assuming that musical features can be properly mapped to visual contents through semantic links between the two domains, we solve the music-to-visual style transfer problem in two steps: music visualization and style transfer. The music visualization network utilizes an encoder-generator architecture with a conditional generative adversarial network to generate image-based music representations from music data. This network is integrated with an image style transfer method to accomplish the style transfer process. Experiments are conducted on WikiArt-IMSLP, a newly compiled dataset including Western music recordings and paintings listed by decades. By utilizing such a label to learn the semantic connection between paintings and music, we demonstrate that the proposed framework can generate diverse image style representations from a music piece, and these representations can unveil certain art forms of the same era. Subjective testing results also emphasize the role of the era label in improving the perceptual quality on the compatibility between music and visual content.

## CCS CONCEPTS

• **Applied computing** → **Media arts**; • **Human-centered computing** → *Visualization*; • **Computing methodologies** → *Machine learning*.

## KEYWORDS

Deep learning; style transfer; generative adversarial networks; aesthetics assessment

*First two authors contributed equally to this research.

## 1 INTRODUCTION

Since Gatys *et al.* proposed the neural algorithm for image style transfer [11], deep learning-based style transfer has been extensively studied. Various types of neural network models now can

modify the texture of an image [13, 19, 30], the genre or instrument of a music piece [20, 21], and the sentiment of texts [29], with all their content information being preserved. Despite its success, one notable issue in these style transfer methods is that almost all of them operate merely within one single data modality, e.g., from one image to another image, or from one music piece to another. Such a situation is far from the cases that human beings design, create, and interpret an artwork, where good ideas and inspirations usually stem from the interplay among the materials from different data modalities. It is quite natural for human artists to break the restriction of data modality, by projecting their imagination of a music piece into their paintings, or by altering a paragraph of texts in a novel into a scene in a movie. In such processes, one would consider a generation problem mapping from one data modality to another via a latent space, and this latent space properly encodes the styles shared by both sides. This problem is referred to as *cross-modal style transfer* in this paper.

Previous investigation of cross-modal learning has been mostly focused on *content generation* rather than style transfer. Most of these studies leverage the techniques of deep transfer learning for various tasks, such as generating images from sound or sounds from images [33, 36, 37]. In comparison to cross-modal content generation, endeavors to cross-modal style transfer are still rarely investigated. However, cross-modal style transfer is often a critical part in the practice of creativity. For example, in design products of virtual reality, animation, and interactive arts, to synergize the styles of visual and music contents is a complicated job, and an automatic process could greatly reduce the efforts. Cross-modal style transfer therefore opens a much broader yet unexplored field for deep generation models.

In this paper, we for the first time investigate *music-to-visual style transfer*, a cross-modal style transfer task considered as image style transfer with music as extra condition. This task is conceptually demonstrated in Figure 1, where an Impressionism music piece is encoded to modify the color scheme an American painting in the mid-19th century into an Impressionism-like one. Such utility demonstrates great potential in integrating visual and music contents in animation, virtual reality, and real-world environments such as concerts.

The music-to-visual style transfer network is proposed based on an assumption that music and image styles can be properly linked by their shared semantic labels. With this assumption, three major questions need to be answered: 1) how to link such semantic labels together, 2) how to evaluate the efficacy of each network component in achieving such aesthetic quality by learning such semantic labels, and 3) how to evaluate the aesthetic quality of the results. We will answer the first question in Section 3 and 4, by introducing the proposed dataset and the music visualization network. Questions 2 and 3 will be discussed in Section 5.

## 2 RELATED WORKS

### 2.1 Visual style transfer

Given a content image and a style image as input, an image style transfer network typically incorporates two tasks: reconstructing the content image or its representation and approximating the statistics of the texture representations (e.g., the Gram matrix) of

the style image [11]. Notable developments include the Gram matrix of feature maps [11], adaptive instance normalization (AdaIN) [13], whitening and coloring transform (WCT) [19], and patch-based methods [30]. While early developments in style transfer were usually limited by the speed of inference and the classes of output style, recent state-of-the-art style transfer methods have overcome these issues and have shown great potential in online video artistic style transfer [7, 10].

### 2.2 Music style transfer

It is hard to give a holistic definition of music style. Music styles depend on the semantic domain being discussed, such as timbre, performance, or composition styles. Timbre style transfer is usually audio-to-audio style transfer which aims at modifying the timbre such as instrument [9, 20, 32] or the gender of singers' voice [17, 34]. Performance style transfer can be either audio-to-audio or symbolic-to-audio, the latter such as piano performance rendering [12, 22] refers to the tasks of converting deadpan performance data (e.g., MIDI) into expressive performance with a specific interpretation of timing and dynamics. Finally, composition style transfer is usually a symbolic-to-symbolic style transfer problem, which aims at modifying the harmonic, rhythmic or structural attributes of music at the score level, and is applied in music genre transfer [5, 21, 23] or blending [15].

### 2.3 Audio-visual content generation and style transfer

Deep learning has provided unprecedented flexibility in various cross-modal content generation tasks operating among audio, visual, and textual information. Both audio-to-visual or visual-to-audio content generation has been studied. [36] proposed source separation based on the visual information in music performance. [33] proposed an image generation framework taking sound as input. [37] proposed to generate ambient sound or soundscape from a given image.

Cross-model style transfer is still a rarely investigated topic by now. Recently, [6] proposed a text-to-image style transfer framework. Some performance generation works can be regarded as text-to-music style transfer, and one important aim of performance generation is to add expressiveness to a dead-pan performance while preserving its content. How artists and composers transfer visual contents into musical ideas and vice versa has long been an attractive topic in art studies [4, 16]. Continuous efforts for centuries to unravel the relationship between music and visual contents have also engendered new art forms and tools, such as color music [26], Lumia arts [8] and *music visualization* techniques which are widely seen in a modern multimedia world [3, 24].

## 3 DATASET

To make a machine learn the relationship between the styles of visual arts and music, we need a dataset containing music and images sharing the same set of semantic labels. How to present such connection between music and images is challenging. The most straightforward way to achieve this goal might be taking the era (i.e., the years the music was composed or the visual artwork was made) as the shared label, as the first attempt to achieve this
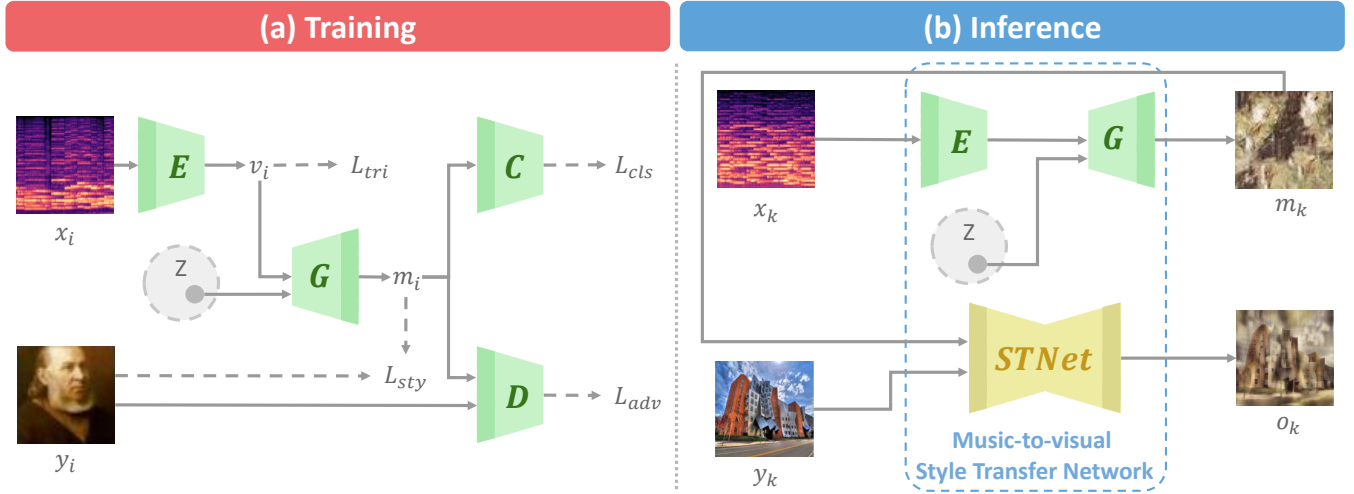
**Figure 2: Overview of the music-to-visual style transfer framework. (a) The training scheme of the music visualization net (MVNet), which contains the encoder-generator pair $\{E, G\}$. The classifier $C$ and discriminator $D$ are for regularization and adversarial training. (b) The inference scheme. The trained MVNet can be integrated with an arbitrary style transfer network for image style transfer, based on the visualized music representation outputted by the MVNet.**

challenge.[1] Other kinds of labels such as genres, user preference, and emotions are likely to limit the size of the dataset because of the incompatibility of labels in different domains. We consider compiling our dataset from two resources: the WikiArt archive[2] for Western visual arts, and the International Music Score Library Project (IMSLP) online music library[3] for Western classical music, both of which contain era label. For the paintings, we use the `Wikiart Retriever`[4] to obtain the images and the corresponding meta-data. As for the music, we use `Selenium`[5] to retrieve classical music pieces in the IMSLP music library. We choose the data from 1480 to the present and annotate their era labels by decades. For example, the music and paintings in 1700-1710 share one label, and those in 1710-1720 share another label. The proposed dataset, named as the WikiArt-IMSLP dataset hereafter, contains in total 11,127 music pieces and 62,968 paintings divided into 54 classes. To the best of our knowledge, this dataset is the first open-source dataset which pairs the music and visual art together. The dataset will be released after the paper is accepted.

Note that the labels in the WikiArt-IMSLP dataset are imbalance. To facilitate the training process, we choose only portraits (which is the largest category in the WikiArt archive) for training. For music, we choose up to 100 music pieces in each era. As a result, there are 11,078 images and 5,587 music pieces for training.

## 4 MUSIC-TO-VISUAL STYLE TRANSFER

We solve the music-to-visual style transfer problem with two steps, namely music visualization and style transfer. Figure 2 illustrates

the pipeline of the proposed system. The system contains two major networks, which are referred to as the *music visualization net* (MVNet) and the *style transfer net* (STNet) in this paper. The MVNet is a regularized encoder-decoder network; its input is an audio data representation, and its output is an image which resembles the style of that image paired with the audio. This image will be referred to as the *style image* hereafter. The style image generated by the MVNet and the target image (i.e. content image) are then fed into the STNet. The output of the STNet is a modified image which resembles the style of the style image.

In what follows, we consider the training data with $N$ music-image pairs $\{x_i, y_i\}_{i=1}^N$, where the $x$ being the 2-D mel-spectrogram of music signals and $y$ being the corresponding image, such that each $x_i$ and $y_i$ were created in the same era. For simplicity, the dimension of all the images is adjusted to $64 \times 64$. The music signals are clips with the length of 8.91 seconds, segmented at the first one-third of each music piece in the dataset. The sampling rate of the music signals is 22.05 kHz. Hamming window with the size of 1024 and hop size of 256 are used for computing spectrogram. The size of the mel-filterbank is 128. The mel-spectrogram is divided into three parts, each of which with 2.97 seconds length. The three parts are then assigned to the three input channels, resulting in the dimension of $128 \times 256 \times 3$. The mel-spectrogram is obtained with the `librosa` library.[6]

The main reason that we divide a mel-spctrogram into three channels for training rather than using merely one channel is described as follows. Our pilot study showed that a single-channel feature is insufficient in representing music information possibly because it is too short. Considering a longer piece of music could solve this issue, but in this case, the input dimension is too large and the model becomes unstable to converge. One compromise is

---

Table 1: The architecture of the networks adopted in this work.

| $x_i \in \mathbb{R}^{128 \times 256} \sim X$ | |
| --- | --- |
| Conv3 × 3-BN-ReLU | 128 × 256 × 32 |
| MaxPool-Conv3 × 3-BN-ReLU | 64 × 128 × 64 |
| MaxPool-Conv3 × 3-BN-ReLU | 32 × 64 × 128 |
| MaxPool-Conv3 × 3-BN-ReLU | 16 × 32 × 256 |
| MaxPool-Conv3 × 3-BN-ReLU | 8 × 16 × 256 |
| Conv1 × 1 | 8 × 16 × 1 |

(a) Encoder

| $\{v_i \oplus z\} \in \mathbb{R}^{8 \times 8 \times 257}$ | |
| --- | --- |
| DeConv4 × 4-IN-ReLU | 16 × 16 × 256 |
| DeConv4 × 4-IN-ReLU | 32 × 32 × 128 |
| Self attention module | 32 × 32 × 128 |
| Self attention module | 32 × 32 × 128 |
| DeConv4 × 4-Tanh | 64 × 64 × 3 |

(b) Generator

| $y_i \in \mathbb{R}^{64 \times 64 \times 3} \sim Y$ | |
| --- | --- |
| Conv4 × 4-LeakyReLU | 32 × 32 × 64 |
| Conv4 × 4-LeakyReLU | 16 × 16 × 128 |
| Conv4 × 4-LeakyReLU | 8 × 8 × 256 |
| Conv4 × 4-LeakyReLU | 4 × 4 × 256 |

(c) Discriminator

to use a long music segment but split it into three channels. In this case, the role of the convolutional filters in the network would differ from the usual ones in processing the RGB channels in images, as now the mel-spectrograms in the three channels are not necessarily synchronized.

## 4.1 Training the Music Visualization Net (MVNet)

The left part of Fig. 2 demonstrates the MVNet in the training stage, which contains an encoder $E$, a generator $G$, a discriminator $D$, and an auxiliary classifier $C$. The encoder $E$ first encodes $x_i$ into a latent vector $v_i := E(x_i)$. A triplet loss term $L_{tri}$ [28] regularizes the behaviors of $v_i$ such that any two $v_i$s encoded from the music pieces in the same era are as similar to each other as possible, while any two $v_i$s from different era are as far to each other as possible. That means, for every input triplet $(x_i^a, x_i^p, x_i^n)$ where $x_i^p$ and $x_i^a$ are in the same era, $x_i^n$ and $x_i^a$ are in different eras, and $(v_i^a, v_i^p, v_i^n) = (E(x_i^a), E(x_i^p), E(x_i^n))$, the triplet loss is represented as

$$L_{tri} = \sum_{i=1}^{N} \left[ \|v_i^a - v_i^p\|_2^2 - \|v_i^a - v_i^n\|_2^2 + \alpha \right]_+ \quad (1)$$

where $\alpha$ is the margin which is set to 1.0 in our experiment, and $[\cdot]_+ := \max(0, \cdot)$. In the training process, we select $x_i^p$ and $x_i^n$ by randomly sampling images from training set.

The latent vector $v_i$ is then concatenated with a random vector $z \sim \mathcal{N}(0, I)$ and fed into the generator $G$. The output $m_i$ is represented $m_i = G(v_i \oplus z)$. To make $m_i$ resemble the style of $y_i$, we impose a style loss term $L_{sty}$ on the generator network. The style loss is defined as the $l_1$ distance in Gram matrices between two images, and we follow the method in [11] to compute the style loss based on the feature maps of a pre-trained VGG net [31]:

$$L_{sty} = \sum_{i=1}^{N} \sum_{s=1}^{S} \|Gram(VGG_s(m_i)) - Gram(VGG_s(y_i))\|_1 \quad (2)$$

where $S$ is the number of layers, $Gram(\cdot)$ is the Gram matrix operation [11], and $VGG_s(\cdot)$ represents the $s$th-layered feature map of the pre-trained VGG net.

We also introduce an adversarial loss $L_{adv}$ to enhance the training process. Following [20], we employ RaGAN [14] as our adversarial training mechanism. Let $\bar{D}$ represent the discriminator $D$ without the last sigmoid layer; that means, $D(a, b) := \text{sigmoid}(\bar{D}(a) - \bar{D}(b))$. Let $P_{real}$ and $P_{gen}$ be the distributions of the image and the music representation, respectively. The adversarial loss is represented as

$$L_{adv} = -\mathbb{E}_{y_i \sim P_{real}}[\log(\mathbb{E}_{m_i \sim P_{gen}}[D(y_i, m_i)])]$$
$$- \mathbb{E}_{m_i \sim P_{gen}}[\log(1 - \mathbb{E}_{y_i \sim P_{real}}[D(m_i, y_i)])] \quad (3)$$

Finally, to better utilize the era labels in the training data, we further introduce a classifier $C$ for era classification. By using this classifier we expect that the likelihood function of the music representation $P(c|m_i)$ should approximate the likelihood function of the training image $P(c|y_i)$, where $c$ represents the era labels. We therefore consider the era classification loss $L_{cls}$ for $y_i$ and $m_i$:

$$L_{cls}^C = \sum_{i=1}^{N} \log P(c|y_i) \quad (4)$$

$$L_{cls}^G = \sum_{i=1}^{N} \log P(c|m_i) \quad (5)$$

where the superscripts $C$ and $G$ indicate the sub-network being updated when that loss term is used during training. More specifically, when training the auxiliary classifier $C$, we adopt $L_{cls}^C$ in Equation 4 to fit the training image $y_i$ their labels; when training $G$, we adopt $L_{cls}^G$ in Equation 5 to fit the generated music representation $m_i$ to the label of their corresponding music piece $x_i$. In this way, the generator is regularized by this loss term since $x_i$ and $y_i$ are paired with the same set of era labels.

In summary, the total loss function $L$ for training the network $\{E, G, D\}$ is represent as

$$L = L_{adv} + \alpha L_{tri} + \beta L_{sty} + \gamma L_{cls}^G, \quad (6)$$

and the era classifier $C$ is trained solely with $L_{cls}^C$. In our experiment, we set $\alpha = 0.1$, $\beta = 10.0$ and $\gamma = 0.1$. In the following, we refer to the trained encoder-decoder pair $\{E, G\}$ as the MVNet.

Note that the training pair $\{x_i, y_i\}$ is not uniquely defined. One $x_i$ can be paired with multiple $y_i$ having the same era label to $x_i$. Therefore, to include more possible $\{x_i, y_i\}$ pairs, these pairs are randomly shuffled for every training epoch.

**Charles Wesley,**
**String Quartets,**
**(1778)**

**Joachim Raff,**
**String Quartet No.5, Op.138,**
**(1867)**

**Eugene Goossens,**
**String Quartet No.2, Op.59,**
**(1942)**

**Christian Junck,**
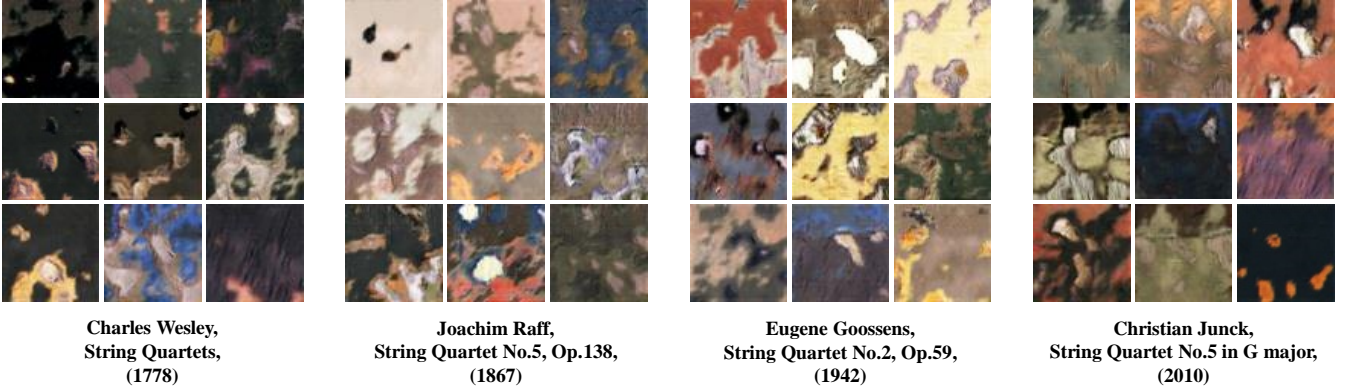**String Quartet No.5 in G major,**
**(2010)**

Figure 3: Random music representation samples generated from four string quartets composed in different eras. The composers' name, the title of music, and the year of composition are listed below the generated samples.

Table 1 shows the components of the encoder, generator and discriminator. Each row of the tables lists the operations employed in each layer (left, different operators are separated with hyphen), and the size of output feature (right). The kernel size for each max pooling operation is two. The structure of the encoder is similar to the encoder part of a traditional U-Net [27]. The architecture of generator generally follows the self-attention GAN (SAGAN) [35], where we employ the self-attention module after the stacked deconvolution layers. For each deconvolution operator in the generator, we adopt spectral normalization to guarantee the training stability. It should be noted that the music latent vector $v_i$ is resized to $8 \times 8 \times 1$ with bi-linear interpolation, in order to guarantee the consistency among feature dimensions, since the size of the random vector $z$ is $8 \times 8 \times 256$. As a result, the overall size of the feature $\{v_i \oplus z\}$ is $8 \times 8 \times 257$. At last, to enforce better modeling of high-frequency structure, we follow the idea of PatchGAN [38], and have the discriminator output a feature map with size of $4 \times 4 \times 256$.

### 4.2 Inference

The right part of Figure 2 demonstrates the inference procedure. After finishing training, the trained MVNet $\{E, G\}$ is employed to generate images (i.e. music representation $m_i$) from music, and the images are expected to convey visual style information of the paintings coming from the era of that music piece.

This style image and the content image are fed into the STNet, which can be an arbitrary image style transfer network, such as those being reviewed in the previous sections. More specifically, for a given content image $y_k$ and an arbitrary music segment $x_k$, the stylized image $o_k$ can be generated as:

$$o_k = STNet(y_k, G(E(x_k))). \tag{7}$$

A detailed comparison of various STNets can be found in the next section.

## 5 EXPERIMENT AND RESULT

The model is trained using two GTX-1080 Ti GPUs and 2 TB SSD. Training the MVNet takes around 48 hours to acheive convergence.

The system is implemented with Python3.6 and PyTorch deep learning framework. For network optimization, we use the Adam optimizer ($\beta_1 = 0.9$, $\beta_2 = 0.999$), a fixed learning rate 0.0001, and a batch size 16 to train the network. Weight decay is set to 0.0001 to avoid over-fitting. Before the painting images are fed into the network, we resize the image as $96 \times 96$, and randomly crop it with a $64 \times 64$ patch.

Experimental results are demonstrated as follows: First, generated samples of the music representation are illustrated. Second, the performance of the system is assessed through the accuracy of era classification, Third, transferred image and audio samples conditioned on music in different eras and processed with different STNets are illustrated and discussed. Finally, we conduct a systematic user study to evaluate the aesthetic quality of the music-to-visual style transfer system. The supplementary materials, demo images and videos, data, and codes are available at the project webpage: https://sunnerli.github.io/Cross-you-in-style/

### 5.1 Illustration of music representations

To see how the music representations look like, in Figure 3 we illustrate four sets of music representations, which are generated from four string quartets composed in the years of 1778, 1867, 1942, and 2010, respectively. First, we investigate whether the MVNet can generate diverse outputs by random sampling over $z \sim \mathcal{N}(0, I)$. To verify this, nine samples are selected for each set. The illustrated samples indicate that the MVNet does generate diverse outputs, as none of them look the same as others.

Second, by fixing the inputs to be string quartets, Figure 3 allows us to compare how the *styles* in both music and paintings affect each other through neural network mapping. The compositional styles of the four music pieces are quite different: Wesley is in the Classical period; Raff's work is mixed with the Romantic penchant where transposition and chromatic scales are more actively used; and Goosens and Junck are more or less influenced by post-tonal music and electronic music. These difference can be observed from their mel-spectrograms. The styles of paintings in the four periods are also different; for example, before the mid-19th century, low color saturation and unified color tone are preferred, while after the
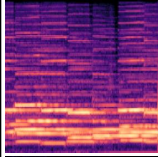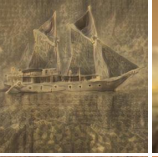
| Composer / title / year | Mel-spectrogram | Music representation | STNet | | | |
| --- | --- | --- | --- | --- | --- | --- |
| | | | AdaIN | WCT | Linear | Avatar-Net |
| St. Jean de Brébeuf, *The Huron Carol* (1643) |  |  |  |  |  |  |
| Johann Caspar Ferdinand Fischer, Ricercar pro Festis Natalytis (1702) |  |  |  |  |  |  |
| Johannes Brahms, Piano Trio No. 1 in B major, Op. 8 (1854) |  |  |  |  |  |  |
| Nathan Shirley, Images (2002) |  |  |  |  |  |  |

Figure 4: Comparison of music-to-image style transfer results over four music pieces from different eras and four STNets. The original content image can be seen in the left of Figure 1. The first column shows the name of the composer, the title of music, and the year of composition. The mel-spectrograms of the music are also illustrated for reference. More results with different content images are provided in supplementary material.



Figure 5: The accuracy of era classification for every epoch. Green line: with classification loss. Orange line: without era classification loss.

mid-19th century, colorful elements (e.g., orange, pink, etc.) with high saturation and complementary colors (e.g., orange and blue, yellow and purple) are characteristic.

We observe that what Figure 3 shows is consistent with the aforementioned historical statements on music and arts. The images of the first set (Wesley's string quartet, 1778) are more similar to each other than the other three sets in terms of tone and saturation. The use of complementary colors within one image is also rarely seen in the first set but commonly seen in the other thr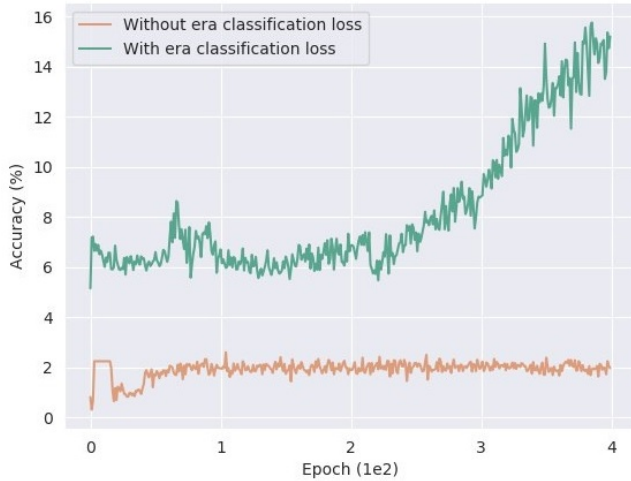ee. This demonstrates that the network does capture the music and image styles and map them to each other properly. More examples of the music representation can be found in Figure 4 and the project webpage.[7]

## 5.2 Era classification

Figure 5 shows the accuracy on era classification computed over the generated music representations for every epoch. Since there are 54 era classes, the accuracy of a random guess is around 2%. Two settings are compared in this experiment: the first setting imposes the classifier loss ($\mathcal{L}_{cls}^{G}$) when training $G$, and the second setting does not include this loss term. The result shows that if the classification loss is not imposed, the accuracy remains at the random guessing level over the epochs. When the classification loss function is imposed, the accuracy increases over the epochs. An increased accuracy implies a higher probability to generate images that can be classified to the correct era through this classifier.

---

[7]The link of project page: https://sunnerli.github.io/Cross-you-in-style/

**Figure 6: Results of music-to-video transfer. The original content images and the transferred results of three selected frames for each video are shown. Music genres from top to down: piano solo, symphony, chamber music, and progressive house. From left to right: composer/ music title/ years, and three pairs of original content image with transferred result which using the linear transformation method.**

## 5.3 Comparison of STNets

To demonstrate the compatibility of the proposed framework with various style transfer methods, we compare the style transfer outputs generated by four different STNets. The STNets include AdaIN [13], WCT [19], linear transformation [18] and Avatar-Net [30]. Figure 4 shows the input mel-spectrograms, the music representations, and the generated results of the four STNets conditioned on four music pieces in different eras. Results show that, again, the generated music representations do correspond with the painting styles in that era: low saturation and unified hue before the Classic period, while high saturation and complementary colors after Romanticism.

For the transferred results, we found that AdaIN merely captures the colors in the music representations, though its processing speed is the fastest among all. The other three STNets can better capture the color scheme of the music representations. The linear transformation can even capture the brushstroke-like texture. WCT tends to emphasize the boundary on a small scale, and Avatar-Net tends to blur the content image. In summary, we indicate that the linear transformation is a compromise between speed and quality. As a result, we will use the linear transformation method in the remaining experiments.

## 5.4 Music-to-video style transfer

In the above discussion, a music segment is assumed to be a static object. We then consider a more realistic case that how music, as an art of time, modifies the visual styles of video in a dynamical manner. As a proof-of-concept study, we consider a preliminary scenario: we select movie clips with background music arrangement, resample the clips to 20 fps, and then for every video frame we take the background music segment around the video frame as the music input to transfer the visual style of that video frame. The length of

the background music segment is 8.91 secs and the middle of this segment is at the time of the video frame. That means, style transfer is processed frame by frame, and each frame takes different but overlapped music segments to generate the music representations. For simplicity, the latent vector $z$ is kept the same over time. We selected four movie clips from *Her* (2013), *Big Fish* (2003), *Scent of a Woman* (1992), and *Trainspotting* (1996) whose background music are piano solo, symphony, chamber music, and progressive house, respectively. The last one does not belong to the classic repertoire. The titles of the music and the composers' name are listed in Figure 6. The movie clips were retrieved from YouTube, and we retrieved these movie clips for research use only.

Figure 6 shows the selected results of music-to-video style transfer. We observe that different music genres transfer the clips into different color tones. Piano solo and symphony transfer the videos into a brighter tone, while chamber music transfer the videos into a darker tone. The genre of progressive house music is obviously not seen in the training set, so the transferred color tone is less usual compared to other samples. A common issue in this task is the fluctuation in brightness and color, which implies that the model is still unstable with the change of music features. To overcome this issue, additional constraints are needed to smooth the style in time domain, and this will be part of our future work.

## 5.5 Aesthetic quality assessment: a user study

To examine the visual aesthetics of the transferred images and videos, we refer to prior research on measuring people's perceived aesthetic impressions [1, 2, 25]. Key scholars summarize various dimensions of aesthetic judgements, from a general dimension "pleasingness" to commonly emphasized cognitive (e.g., comprehensibility, originality) and emotional (e.g., emotiveness, impressiveness) dimensions. We adopted questions from the
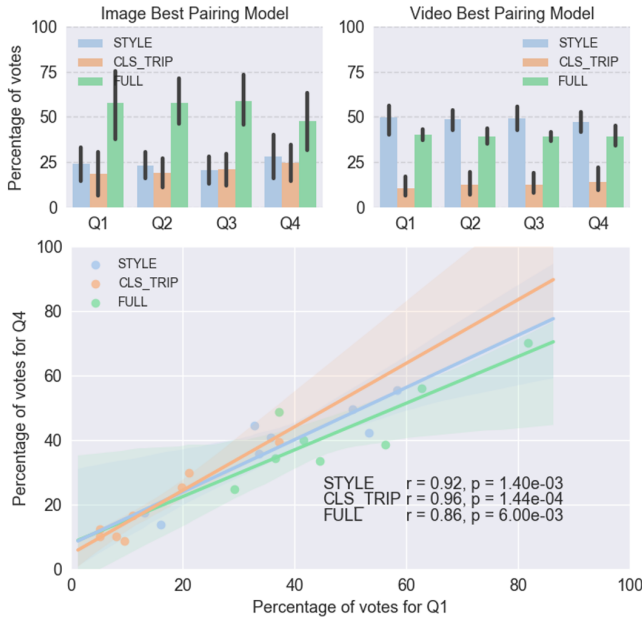
**Figure 7: User study result. Upper left: best pairing model (in %) on M2I samples. Upper right: best pairing model (in %) on M2V samples. Bottom: correlation between Q1 and Q4.**

emotional and pleasingness dimensions as the first step to assess visual aesthetics of the images (including painting, photograph) and videos after the transfer. We decided to focus on general perceived aesthetic impression instead of more technical aspects of visual aesthetics (e.g., saturation, contrast, stroke), as the latter type of assessment requires knowledge at the expert level.

We recruited a total of 137 participants (47.20% male, 52.10% female) to fill out an online questionnaire to assess eight image/video samples. Four are music-to-image (M2I) while four are music-to-video (M2V) samples. First, we had participants get familiar with the original images and audio/video content. Then, participants were asked to compare the style transfer outputs generated from three variants of our model: 1) style loss only (STYLE), 2) style loss plus the triplet and classification loss (CLS_TRIP), and 3) the three loss terms plus the adversarial loss (FULL). It should be noted that we compared the original images/videos and STYLE in a pilot study but found that the clean images/videos could bias human judgement on the artifacts of the style transfer results. Therefore, we removed the images/ videos in the official study to reduce survey completion difficulty. For each image/ audio/ video content, participants were asked to select the one they perceived the most beautiful (Q1), attractive (Q2), moving (Q3), and most harmonious with the background music (Q4) out of the three images/videos style transfer outputs. Q1 and Q2 belong to the pleasingness dimension; Q3 belongs to the emotion dimension; Q4 is the objective of our research.

Results are shown in Figure 7. For M2I samples, participants preferred the FULL model the most while for M2V samples, participants preferred the STYLE model. The difference in M2I and M2V reveals the roles of the objective functions in the generation

process. Using only style loss in the M2I case seems to generate âĂŸboringâĂŹ result, but in the M2V case this is preferred as adding additional loss terms usually results in fluctuations in video. Comparing FULL to CLS_TRIP, it shows that adding the adversarial loss helps both M2I and M2V. Additionally, participant responses to the four questions strongly correlated to each other. For instance, the correlation between the percentage of votes of Q1 and Q4 suggests that the model making the background music more harmonious with the visual style tends to make such visual style more beautiful (r = 0.92, 0.96 and 0.86 for the three models, respectively).

## 6 CONCLUSION

We have demonstrated the feasibility of transferring visual styles directly from music. The flexibility of using different style transfer networks and of using either image or video contents in our framework all suggest great potential in the applications of animation and interactive arts. Evaluation results indicate the importance of a shared semantic space in solving the cross-modal style transfer problem, and also reveal the multi-dimensional nature of aesthetic quality assessment, which is still a challenging problem worth further study.

We have emphasized the importance of cross-model transfer in human creativity process. However, it should be noted that our proposed solution does not include all the scenarios that human artists deal with this problem in the real world. In fact, shared labels are not a necessary condition in real-world cross-model style transfer process. The condition of pairing an image to a music piece can usually be arbitrary and relies on how one interprets it. Shared labels are neither a sufficient condition in real-world cross-model style transfer. The reason that we can imagine a visual scene when listening to music is that our brains have built a complicated web of meaning that connects these data in various semantic levels, not because we know the time they were composed or painted. Annotations of high-level semantics such as art movement or genre could partly address this issue, but such annotations might be more difficult to be shared. On the other hand, the era labels adopted in this work ignore the the time asynchronization in the development of art and music (e.g., impressionism music appeared later than impressionism art). The purpose of this work is not to claim that learning the semantic links from era labels is the unique and 'correctâĂŹ way to assign a style to an artwork. Rather, we emphasize that learning the arbitrary semantic links between different domains can be a feasible and scalable way for content generation. A future direction toward a more advanced cross-modal style transfer is to establish more label information such as genres and emotions to link the data from different modalities altogether up to a higher semantic level.

# REFERENCES

[1] M Augustin, Claus-Christian Carbon, and Johan Wagemans. 2011. Measuring aesthetic impressions of visual art. *PERCEPTION* 40 (01 2011), 219.

[2] M Dorothee Augustin, Johan Wagemans, and Claus-Christian Carbon. 2012. All is beautiful? Generality vs. specificity of word usage in visual aesthetics. *Acta Psychologica* 139, 1 (2012), 187–201.

[3] Tony Bergstrom, Karrie Karahalios, and John C Hart. 2007. Isochords: visualizing structure in music. In *ACM Proc. Graphics Interface*. 297–304.

[4] Greta Berman. 1999. Synesthesia and the Arts. *Leonardo* 32, 1 (1999), 15–22.

[5] Gino Brunner, Yuyi Wang, Roger Wattenhofer, and Sumu Zhao. 2018. Symbolic music genre transfer with cyclegan. In *ICTAI*. 786–793.

[6] Sahil Chelaramani, Abhishek Jha, and Anoop M. Namboodiri. 2018. Cross-Modal Style Transfer. In *ICIP*. 2157–2161.

[7] Dongdong Chen, Jing Liao, Lu Yuan, Nenghai Yu, and Gang Hua. 2017. Coherent online video style transfer. In *ICCV*. 1105–1114.

[8] Fred Collopy. 2000. Color, form, and motion: Dimensions of a musical art of light. *Leonardo* 33, 5 (2000), 355–360.

[9] Jesse Engel, Cinjon Resnick, Adam Roberts, Sander Dieleman, Mohammad Norouzi, Douglas Eck, and Karen Simonyan. 2017. Neural audio synthesis of musical notes with wavenet autoencoders. In *ICML*. 1068–1077.

[10] Chang Gao, Derun Gu, Fangjun Zhang, and Yizhou Yu. 2018. ReCoNet: Real-time Coherent Video Style Transfer Network. In *ACCV*. 637–653.

[11] Leon A Gatys, Alexander S Ecker, and Matthias Bethge. 2016. Image style transfer using convolutional neural networks. In *CVPR*. 2414–2423.

[12] Curtis Hawthorne, Anna Huang, Daphne Ippolito, and Douglas Eck. 2018. Transformer-NADE for Piano Performances. In *NIPS 2nd Workshop on Machine Learning for Creativity and Design*.

[13] Xun Huang and Serge Belongie. 2017. Arbitrary style transfer in real-time with adaptive instance normalization. In *ICCV*. 1501–1510.

[14] Alexia Jolicoeur-Martineau. 2018. The relativistic discriminator: a key element missing from standard GAN. *arXiv preprint arXiv:1807.00734* (2018).

[15] Maximos Kaliakatsos-Papakastas, Marcelo Queiroz, Costas Tsougras, and Emilios Cambouropoulos. 2017. Conceptual blending of harmonic spaces for creative melodic harmonisation. *JNMR* 46, 4 (2017), 305–328.

[16] Sharon L Kennedy. 2007. Painting music: rhythm and movement in art. (2007).

[17] Kazuhiro Kobayashi, Tomoki Toda, Graham Neubig, Sakriani Sakti, and Satoshi Nakamura. 2014. Statistical singing voice conversion with direct waveform modification based on the spectrum differential. In *ISCA*.

[18] Xueting Li, Sifei Liu, Jan Kautz, and Ming-Hsuan Yang. 2018. Learning linear transformations for fast arbitrary style transfer. *arXiv:1808.04537* (2018).

[19] Yijun Li, Chen Fang, Jimei Yang, Zhaowen Wang, Xin Lu, and Ming-Hsuan Yang. 2017. Universal style transfer via feature transforms. In *NeurIPS*. 386–396.

[20] Chien-Yu Lu, Min-Xin Xue, Chia-Che Chang, Che-Rung Lee, and Li Su. 2018. Play as You Like: Timbre-enhanced Multi-modal Music Style Transfer. *arXiv preprint arXiv:1811.12214* (2018).

[21] Wei Tsung Lu and Li Su. 2018. Transferring the Style of Homophonic Music Using Recurrent Neural Networks and Autoregressive Model.. In *ISMIR*. 740–746.

[22] A Maezawa. 2018. Deep piano performance rendering with conditional VAE. In *ISMIR Late Breaking and Demo Papers*.

[23] Iman Malik and Carl Henrik Ek. 2017. Neural translation of musical style. *arXiv preprint arXiv:1708.03535* (2017).

[24] Arpi Mardirossian and Elaine Chew. 2007. Visualizing Music: Tonal Progressions and Distributions.. In *ISMIR*. Citeseer, 189–194.

[25] Morten Moshagen and Meinald T. Thielsch. 2010. Facets of Visual Aesthetics. *Int. J. Hum.-Comput. Stud.* 68, 10 (Oct. 2010), 689âĂŞ709. https://doi.org/10.1016/j.ijhcs.2010.05.006

[26] Kenneth Peacock. 1988. Instruments to perform color-music: Two centuries of technological experimentation. *Leonardo* 21, 4 (1988), 397–406.

[27] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. 2015. U-net: Convolutional networks for biomedical image segmentation. In *International Conference on Medical image computing and computer-assisted intervention*. Springer, 234–241.

[28] F. Schroff, D. Kalenichenko, and J. Philbin. 2015. FaceNet: A unified embedding for face recognition and clustering. In *CVPR*. 815–823. https://doi.org/10.1109/CVPR.2015.7298682

[29] Tianxiao Shen, Tao Lei, Regina Barzilay, and Tommi Jaakkola. 2017. Style transfer from non-parallel text by cross-alignment. In *NeurIPS*. 6830–6841.

[30] Lu Sheng, Ziyi Lin, Jing Shao, and Xiaogang Wang. 2018. Avatar-net: Multi-scale zero-shot style transfer by feature decoration. In *CVPR*. 8242–8250.

[31] Karen Simonyan and Andrew Zisserman. 2014. Very Deep Convolutional Networks for Large-Scale Image Recognition. *arXiv 1409.1556* (09 2014).

[32] Prateek Verma and Julius O Smith. 2018. Neural style transfer for audio spectograms. *arXiv preprint arXiv:1801.01589* (2018).

[33] Chia-Hung Wan, Shun-Po Chuang, and Hung-Yi Lee. 2019. Towards Audio to Scene Image Synthesis Using Generative Adversarial Network. In *ICASSP*. 496–500.

[34] Cheng-Wei Wu, Jen-Yu Liu, Yi-Hsuan Yang, and Jyh-Shing R Jang. 2018. Singing Style Transfer Using Cycle-Consistent Boundary Equilibrium Generative Adversarial Networks. *arXiv preprint arXiv:1807.02254* (2018).

[35] Han Zhang, Ian Goodfellow, Dimitris Metaxas, and Augustus Odena. 2018. Self-attention generative adversarial networks. *arXiv preprint arXiv:1805.08318* (2018).

[36] Hang Zhao, Chuang Gan, Andrew Rouditchenko, Carl Vondrick, Josh McDermott, and Antonio Torralba. 2018. The sound of pixels. In *ECCV*. 570–586.

[37] Yipin Zhou, Zhaowen Wang, Chen Fang, Trung Bui, and Tamara L Berg. 2018. Visual to sound: Generating natural sound for videos in the wild. In *CVPR*. 3550–3558.

[38] Jun-Yan Zhu, Taesung Park, Phillip Isola, and Alexei A Efros. 2017. Unpaired image-to-image translation using cycle-consistent adversarial networks. In *Proceedings of the IEEE international conference on computer vision*. 2223–2232.