

MBZUAI

Digital.Commons@MBZUAI

Computer Vision Faculty Publications

Scholarly Works

4-8-2021

Handwriting Transformers

Ankan Kumar Bhunia

Mohamed bin Zayed University of Artificial Intelligence

Salman Khan

Mohamed bin Zayed University of Artificial Intelligence & Australian National University

Hisham Cholakkal

Mohamed bin Zayed University of Artificial Intelligence

Rao Muhammad Anwer

Mohamed bin Zayed University of Artificial Intelligence

Fahad Shahbaz Khan

Mohamed bin Zayed University of Artificial Intelligence & Linköping University

See next page for additional authors

Follow this and additional works at: <https://dclibrary.mbzuai.ac.ae/cvfp>



Part of the [Computer Sciences Commons](#)

Preprint: arXiv

- Archived with thanks to arXiv
- Preprint License: [CC by 4.0](#)
- Uploaded 24 March 2022

Recommended Citation

A.K. Bhunia, S. Khan, H. Cholakkal, R.M. Anwer, F.S. Khan, and M.A. Shah, "Handwriting transformers", 2021, arXiv:2104.03964

This Article is brought to you for free and open access by the Scholarly Works at Digital.Commons@MBZUAI. It has been accepted for inclusion in Computer Vision Faculty Publications by an authorized administrator of Digital.Commons@MBZUAI. For more information, please contact libraryservices@mbzuai.ac.ae.

Authors

Ankan Kumar Bhunia, Salman Khan, Hisham Cholakkal, Rao Muhammad Anwer, Fahad Shahbaz Khan, and Mubarak A. Shah

Handwriting Transformers

Ankan Kumar Bhunia¹ Salman Khan^{1,2} Hisham Cholakkal¹ Rao Muhammad Anwer¹
 Fahad Shahbaz Khan^{1,3} Mubarak Shah⁴
¹Mohamed bin Zayed University of AI, UAE ²Australian National University, Australia
³Linköping University, Sweden ⁴University of Central Florida, USA

Abstract

We propose a novel transformer-based styled handwritten text image generation approach, HWT, that strives to learn both style-content entanglement as well as global and local writing style patterns. The proposed HWT captures the long and short range relationships within the style examples through a self-attention mechanism, thereby encoding both global and local style patterns. Further, the proposed transformer-based HWT comprises an encoder-decoder attention that enables style-content entanglement by gathering the style representation of each query character. To the best of our knowledge, we are the first to introduce a transformer-based generative network for styled handwritten text generation.

Our proposed HWT generates realistic styled handwritten text images and significantly outperforms the state-of-the-art demonstrated through extensive qualitative, quantitative and human-based evaluations. The proposed HWT can handle arbitrary length of text and any desired writing style in a few-shot setting. Further, our HWT generalizes well to the challenging scenario where both words and writing style are unseen during training, generating realistic styled handwritten text images.

1. Introduction

Generating realistic synthetic handwritten text images, from typed text, that is versatile in terms of both writing style and lexicon is a challenging problem. Automatic handwritten text generation can be beneficial for people having disabilities or injuries that prevent them from writing, translating a note or a memo from one language to another by adapting an author’s writing style or gathering additional data for training deep learning-based handwritten text recognition models. Here, we investigate the problem of realistic handwritten text generation of unconstrained text sequences with arbitrary length and diverse calligraphic attributes representing writing styles of a writer.

Generative Adversarial Networks (GANs) [8] have been

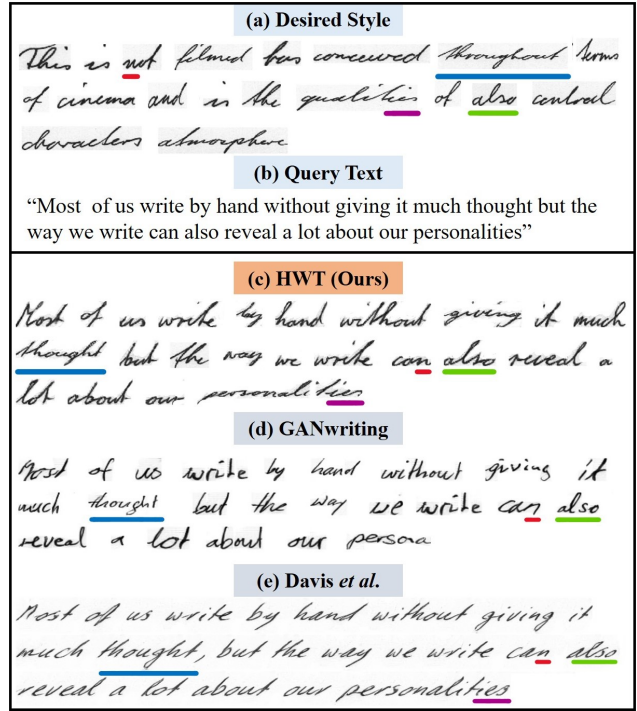


Figure 1: Comparison of HWT (c) with GANwriting [14] (d) and Davis *et al.* [5] (e) in imitating the desired unseen writing style (a) for given query text (b). While [14, 5] capture global writing styles (e.g., slant), they struggle to imitate local style patterns (e.g., character style, ligatures). HWT (c) imitates both global and local styles, leading to a more realistic styled handwritten text image generation. For instance, style of ‘n’ (red line) appearing in (a) is mimicked by HWT, for a different word including same character ‘n’. Similarly, a group of characters in ‘thought’ and ‘personalities’ (blue and magenta lines) are styled in a way that matches with words (‘throughout’ and ‘qualities’) sharing some common characters in (a). Furthermore, HWT preserves cursive patterns and connectivity of all characters in word ‘also’ (green line).

investigated for offline handwritten text image generation [4, 3, 14, 7, 5]. These methods strive to directly synthe-

size text images by using offline handwriting images during training, thereby extracting useful features, such as writing appearance (*e.g.*, ink width, writing slant) and line thickness changes. Alonso *et al.* [3] propose a generative architecture that is conditioned on input content strings, thereby not restricted to a particular pre-defined vocabulary. However, their approach is trained on isolated fixed-sized word images and struggles to produce high quality arbitrarily long text along with suffering from style collapse. Fogel *et al.* [7] introduce a ScrabbleGAN approach, where the generated image width is made proportional to the input text length. ScrabbleGAN is shown to achieve impressive results with respect to the content. However, both [3, 7] do not adapt to a specific author’s writing style.

Recently, GAN-based approaches [5, 14] have been introduced for the problem of styled handwritten text image generation. These methods take into account both content and style, when generating offline handwritten text images. Davis *et al.* [5] propose an approach based on StyleGAN [15] and learn generated handwriting image width based on style and input text. The GANwriting framework [14] conditions handwritten text generation process to both textual content and style features in a few-shot setup.

In this work, we distinguish two key issues that impede the quality of styled handwritten text image generation in the existing GAN-based methods [5, 14]. First, both style and content are loosely connected as their representative features are processed separately and later concatenated. While such a scheme enables entanglement between style and content at the word/line-level, it does not explicitly enforce style-content entanglement at the character-level. Second, although these approaches capture global writing style (*e.g.*, ink width, slant), they do not explicitly encode local style patterns (*e.g.*, character style, ligatures). As a result of these issues, they struggle to accurately imitate local calligraphic style patterns from reference style examples (see Fig. 1). Here, we look into an alternative approach that addresses both these issues in a single generative architecture.

1.1. Contributions

We introduce a new styled handwritten text generation approach built upon transformers, termed Handwriting Transformers (HWT), that comprises an encoder-decoder network. The encoder network utilizes a multi-headed self-attention mechanism to generate a self-attentive style feature sequence of a writer. This feature sequence is then input to the decoder network that consists of multi-headed self- and encoder-decoder attention to generate character-specific style attributes, given a set of query word strings. Consequently, the resulting output is fed to a convolutional decoder to generate final styled handwritten text image. Moreover, we improve the style consistency of the generated text by constraining the decoder output through a

loss term whose objective is to re-generate style feature sequence of a writer at the encoder.

Our HWT imitates the style of a writer for a given query content through self- and encoder-decoder attention that emphasizes relevant self-attentive style features with respect to each character in that query. This enables us to capture style-content entanglement at the character-level. Furthermore, the self-attentive style feature sequence generated by our encoder captures both the global (*e.g.*, ink width, slant) and local styles (*e.g.*, character style, ligatures) of a writer within the feature sequence.

We validate our proposed HWT by conducting extensive qualitative, quantitative and human-based evaluations. In the human-based evaluation, our proposed HWT was preferred 81% of the time over recent styled handwritten text generation methods [5, 14], achieving human plausibility in terms of the writing style mimicry. Following GANwriting [14], we evaluate our HWT on all the four settings on the IAM handwriting dataset. On the extreme setting of out-of-vocabulary and unseen styles (OOV-U), where both query words and writing styles are never seen during training, the proposed HWT outperforms GANwriting [14] with an absolute gain of 16.5 in terms of Fr chet Inception Distance (FID) thereby demonstrating our generalization capabilities. Further, our qualitative analysis suggest that HWT performs favorably against existing works, generating realistic styled handwritten text images (see Fig. 1).

2. Related Work

Deep learning-based handwritten text generation approaches can be roughly divided into stroke-based online and image-based offline methods. Online handwritten text generation methods [9, 2] typically require temporal data acquired from stroke-by-stroke recording of real handwritten examples (vector form) using a digital stylus pen. On the other hand, recent generative offline handwritten text generation methods [4, 3, 14, 7] aim to directly generate text by performing training on offline handwriting images.

Graves [9] proposes an approach based on Recurrent Neural Network (RNN) with Long-Term Memory (LSTM) cells, which enables predicting future stroke points from previous pen positions and an input text. Aksan *et al.* [4] propose a method based on conditional Variational RNN (VRNN), where the input is split into two separate latent variables to represent content and style. However, their approach tends to average out particular styles across writers, thereby reducing details [17]. In a subsequent work [1], the VRNN module is substituted by Stochastic Temporal CNNs which is shown to provide more consistent generation of handwriting. Kotani *et al.* [17] propose an online handwriting stroke representation approach to represent latent style information by encoding writer-, character- and writer-character-specific style changes within an RNN model.

Other than sequential methods, several recent works have investigated offline handwritten text image generation using GANs. Haines *et al.* [11] introduce an approach to generate new text in a distinct style inferred from source images. Their model requires a certain degree of human intervention during character segmentation and is limited to generating characters that are in the source images. The work of [4] utilize CycleGAN [24] to synthesize images of isolated handwritten characters of Chinese language. Alonso *et al.* [3] propose an approach, where handwritten text generation is conditioned by character sequences. However, their approach suffers from style collapse hindering the diversity of synthesized images. Fogel *et al.* [7] propose an approach, called ScrabbleGAN, that synthesizes handwritten word using a fully convolutional architecture. Here, the characters generated have similar receptive field width. A conversion model is introduced by [20] that approximates online handwriting from offline samples followed by using style transfer technique to the online data. This approach relies on conversion model’s performance.

Few recent GAN-based works [5, 14] investigate the problem of offline styled handwritten text image generation. Davis *et al.* [5] propose an approach, where handwritten text generation is conditioned on both text and style, capturing global handwriting style variations. Kang *et al.* [14] propose a method, called GANwriting, that conditions text generation on extracting style features in a few-shot setup and textual content of a predefined fixed length.

Our Approach: Similar to GANwriting [14], we also investigate the problem of styled handwritten text generation in a few-shot setting, where a limited number of style examples are available for each writer. Different from GANwriting, our approach possesses the flexibility to generate styled text of arbitrary length. In addition, existing works [5, 14] only capture style-content entanglement at the word/line-level. In contrast, our transformer-based approach enables style-content entanglement both at the word and character-level. While [5, 14] focuses on capturing the writing style at the global level, the proposed method strives to imitate both global and local writing style.

3. Proposed Approach

Motivation: To motivate our proposed HWT method, we first distinguish two desirable characteristics to be considered when designing an approach for styled handwritten text generation with varying length and any desired style in a few-shot setting, without using character-level annotation. **Style-Content Entanglement:** As discussed earlier, both style and content are loosely connected in recently introduced GAN-based works [14, 5] with separate processing of style and content features, which are later concatenated. Such a scheme does not explicitly encode style-content entanglement at the character-level. Moreover, there are sep-

arate components for style, content modeling followed by a generator for decoding stylized outputs. In addition to style-content entanglement at word/line level, an entanglement between style and content at the character-level is expected to aid in imitating the character-specific writing style along with generalizing to out-of-vocabulary content. Further, such a tight integration between style and content leads to a cohesive architecture design.

Global and Local Style Imitation: While the previous requisite focuses on connecting style and content, the second desirable characteristic aims at modeling both the global as well as local style features for a given calligraphic style. Recent generative methods for styled handwritten text generation [14, 5] typically capture the writing style at the global level (*e.g.*, ink width, slant). However, the local style patterns (*e.g.*, character style, ligatures) are not explicitly taken into account while imitating the style of a given writer. We argue that both global *and* local style patterns are desired to be imitated for accurate styled text image generation.

3.1. Approach Overview

Problem Formulation: We aim to learn the complex handwriting style characteristics of a particular writer $i \in \mathcal{W}$, where \mathcal{W} includes a total of M writers. We are given a set of P handwritten word images, $\mathbf{X}_i^s = \{\mathbf{x}_{ij}\}_{j=1}^P$, as few-shot calligraphic style examples of each writer. The superscript ‘s’ in \mathbf{X}_i^s denotes use of the set as a source of handwriting style which is transferred to the target images $\tilde{\mathbf{X}}_i^t$ with new textual content but consistent style properties. The textual content is represented as a set of input query word strings $\mathcal{A} = \{\mathbf{a}_j\}_{j=1}^Q$, where each word string \mathbf{a}_j comprises an arbitrary number of characters from permitted characters set \mathcal{C} . The set \mathcal{C} includes alphabets, numerical digits and punctuation marks *etc.* Given a query text string $\mathbf{a}_j \in \mathcal{A}$ from an unconstrained set of vocabulary and \mathbf{X}_i^s , our model strives to generate new images $\tilde{\mathbf{X}}_i^t$ with the same text \mathbf{a}_j in the writing style of a desired writer i .

Overall Architecture: Fig. 2 presents an overview of our proposed HWT approach, where a conditional generator G_θ synthesizes handwritten text images, a discriminator D_ψ ensures realistic generation of handwriting styles, a recognizer R_ϕ aids in textual content preservation, and a style classifier S_η ensures satisfactory transfer of the calligraphic styles. The focus of our design is the introduction of a transformer-based generative network for unconstrained styled handwritten text image generation. Our generator G_θ is designed in consideration to the desirable characteristics listed earlier leveraging the impressive learning capabilities of transformer models. To meticulously imitate a handwriting style, a model is desired to learn style-content entanglement as well as global and local style patterns.

To this end, we introduce a transformer-based handwriting generation model, which enables us to capture the long

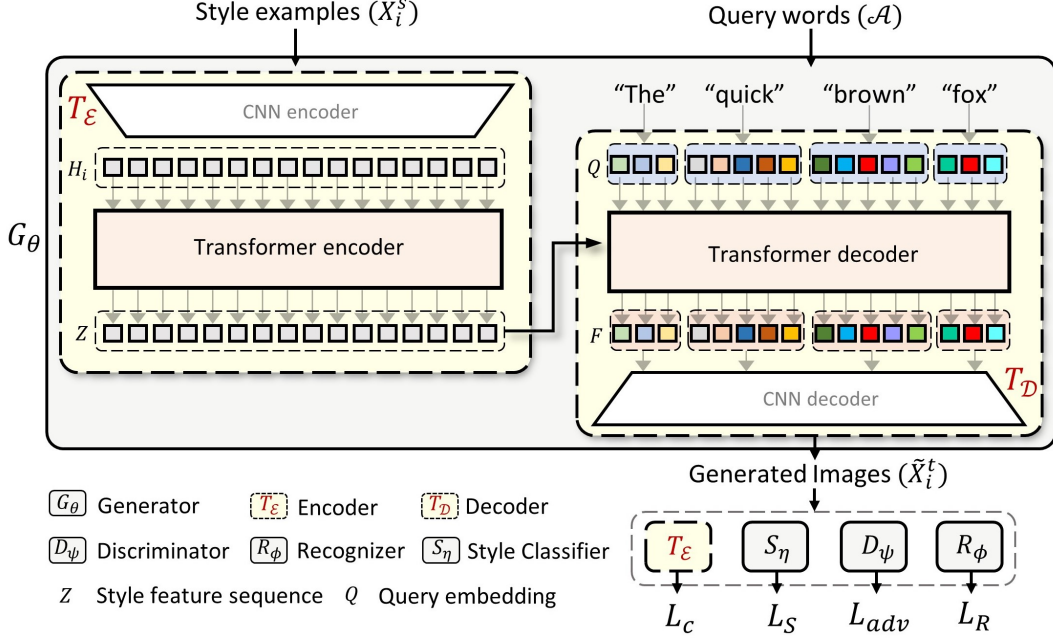


Figure 2: Overall architecture of our Handwriting Transformers (HWT) to generate styled handwritten text images \tilde{X}_i^t . HWT comprises a conditional generator having an encoder T_E and a decoder network T_D . Both the encoder and decoder networks constitute a hybrid convolution and multi-head self-attention design, which combines the strengths of CNN and transformer-based models *i.e.*, highly expressive relationship modeling while working with limited handwriting style example images. Resultantly, our design seamlessly achieves style-content entanglement that encodes relationships between textual content and writer’s style along with learning both global and local style patterns for given inputs (X_i^s and \mathcal{A}).

and short range contextual relationships within the style examples X_i^s by utilizing a self-attention mechanism. In this way, both the global and local style patterns are encoded. Additionally, our transformer-based model comprises an encoder-decoder attention that allows style-content entanglement by inferring the style representation for each query character. A direct applicability of transformer-based design is infeasible in our few-shot setting due to its large data requirements and quadratic complexity. To circumvent this issue, our proposed architecture design utilizes the expressivity of a transformer within the CNN feature space.

The main idea of the proposed HWT method is simple but effective. A transformer-based encoder T_E is first used to model self-attentive style context that is later used by a decoder T_D to generate query text in a specific writer’s style. We define learnable embedding vector $q_c \in \mathbb{R}^{512}$ for each character c of the permissible character set \mathcal{C} . For example, we represent the query word ‘deep’ as a sequence of its respective character embeddings $Q_{\text{deep}} = \{q_d \dots q_p\}$. We refer them as query embeddings. Such a character-wise representation of the query words and the transformer-based sequence processing helps our model to generate handwritten words of variable length, and also qualifies it to produce out-of-vocabulary words more efficiently. Moreover, it avoids averaging out individual character-specific styles in order to maintain the overall (global and local) writing

style. The character-wise style interpolation and transfer is ensured by the self- and encoder-decoder attention in the transformer module that infers the style representation of each query character based on a set of handwritten samples provided as input. We describe the proposed generative architecture in Sec. 3.2 and the loss objectives in Sec. 3.3.

3.2. Generative Network

The generator G_θ includes two main components: an encoder network $T_E : X_i^s \rightarrow Z$ and a decoder network $T_D : (Z, \mathcal{A}) \rightarrow \tilde{X}_i^t$. The encoder produces a sequence of feature embeddings $Z \in \mathbb{R}^{N \times d}$ (termed as style feature sequence) from a given set of style examples X_i^s . The decoder takes Z as an input and converts the input word strings $a_j \in \mathcal{A}$ to realistic handwritten images \tilde{X}_i^t with same style as the given examples X_i^s of a writer i . Both the encoder and decoder networks constitute a *hybrid* design based on convolution and multi-head self-attention networks. This design choice combines the strengths of CNNs and transformer models *i.e.*, highly expressive relationship modeling while working with limited handwriting images. Its worth mentioning that a CNN-only design would struggle to model long-term relations within sequences while an architecture based solely on transformer networks would demand large amount of data and longer training times [16].

Encoder T_E . The encoder aims at modelling both global

and local calligraphic style attributes (*i.e.*, slant, skew, character shapes, ligatures, ink widths *etc.*) from the style examples \mathbf{X}_i^s . Before feeding style images to the highly expressive transformer architecture, we need to represent the style examples as a sequence. A straightforward way would be to flatten the image pixels into a 1D vector [6]. However, given the quadratic complexity of transformer models and their large data requirements, we find this to be infeasible. Instead, we use a CNN backbone network to obtain sequences of convolutional features from the style images. First, we use a ResNet18 [12] model to generate lower-resolution activation maps $\mathbf{h}_{ij} \in \mathbb{R}^{h \times w \times d}$ for each style image \mathbf{x}_{ij} . Then, we flatten the spatial dimension of \mathbf{h}_{ij} to obtain a sequence of feature maps of size $n \times d$, where $n = h \times w$. Each vector in the feature sequence represents a region in the original image and can be considered as the image descriptor for that particular region. After that, we concatenate the feature sequence vectors extracted from all style images together to obtain a single tensor $\mathbf{H}_i \in \mathbb{R}^{N \times d}$, where $N = n \times P$.

The next step includes modeling the global and local compositions between all entities of the obtained feature sequence \mathbf{Z} . A transformer-based encoder is employed for that purpose. The encoder has L layers, where each layer has a standard architecture that consists of a multi-headed self-attention module and a Multi-layer Perceptron (MLP) block. At each layer l , the multi-headed self-attention maps the input sequence from the previous layer \mathbf{H}^{l-1} into a triplet (key \mathbf{K} , query \mathbf{Q} , value \mathbf{V}) of intermediate representations given by,

$$\mathbf{Q} = \mathbf{H}^{l-1} \mathbf{W}^Q, \mathbf{K} = \mathbf{H}^{l-1} \mathbf{W}^K, \mathbf{V} = \mathbf{H}^{l-1} \mathbf{W}^V,$$

where $\mathbf{W}^Q \in \mathbb{R}^{N \times d_q}$, $\mathbf{W}^K \in \mathbb{R}^{N \times d_k}$ and $\mathbf{W}^V \in \mathbb{R}^{N \times d_v}$ are the learnable weight matrix for query, key and value respectively. For each head, the process is represented as,

$$\mathbf{O}^j = \text{softmax} \left(\frac{\mathbf{Q} \mathbf{K}^T}{\sqrt{d_k}} \right) \mathbf{V} \in \mathbb{R}^{N \times d_v}, \quad j \in \{1, \dots, J\}. \quad (1)$$

The concatenation of all J head outputs $\mathbf{O} = [\mathbf{O}^1, \dots, \mathbf{O}^J]$ is then fed through an MLP layer to obtain the output feature sequence \mathbf{H}^l for the layer l . This update procedure is repeated for a total of L layers, resulting in the final feature sequence $\mathbf{Z} \in \mathbb{R}^{N \times d}$. To retain information regarding the order of input sequences being supplied, we add fixed positional encodings [23] to the input of each attention layer.

Decoder T_D . The initial stage in the decoder uses the standard architecture of the transformer that consists of multi-headed self- and encoder-decoder attention mechanisms. Unlike the self-attention, the encoder-decoder attention derives the key and value vectors from the output of the encoder, whereas the query vectors come from the decoder layer itself. For an m_j character word $\mathbf{a}_j \in \mathcal{A}$ (length

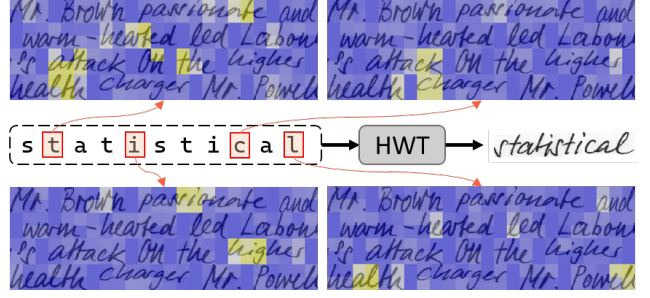


Figure 3: Visualization of encoder-decoder attention maps at the last layer of the transformer decoder. The attention maps are computed for each character in the query word ('statistical') which are then mapped to spatial regions (heat maps) in the example style images. Here, heat maps corresponding to the four different query characters 't', 'i', 'c' and 'l' are shown. For instance, the top-left attention map corresponding to the character 't', highlights multiple image regions containing the character 't'.

m_j being variable depending on the word), the query embedding $\mathbf{Q}_{\mathbf{a}_j} = \{\mathbf{q}_{c_k}\}_{k=1}^{m_j}$ is used as a learnt positional encoding to each attention layer of the decoder. Intuitively, each query embedding learns to look up regions of interest in the style images to infer the style attributes of all query characters (see Fig. 3). Over multiple consecutive decoding layers, these output embeddings accumulate style information, producing a final output $\mathbf{F}_{\mathbf{a}_j} = \{\mathbf{f}_{c_k}\}_{k=1}^{m_j} \in \mathbb{R}^{m_j \times d}$. We process the entire query embedding in parallel at each decoder layer. We add a randomly sampled noise vector $\mathcal{N}(0, 1)$ to the output $\mathbf{F}_{\mathbf{a}_j}$ in order to model the natural variation of individual handwriting. For an m -character word, we concatenate these m_j embedding vectors and pass them through a linear layer, resulting in an $m_j \times 8192$ matrix. After reshaping it to a dimension of $512 \times 4 \times 4m_j$, we pass it through a CNN decoder having four residual blocks followed by a tanh activation layer to obtain final output images (styled hand written text images).

3.3. Training and Loss Objectives

Our training algorithm follows the traditional GAN paradigm, where a discriminator network D_ψ is employed to tell apart the samples generated from generator G_θ from the real ones. As the generated word images are of varying width, the proposed discriminator D_ψ is also designed to be convolutional in nature. We use the hinge version of the adversarial loss [18] defined as,

$$L_{adv} = \mathbb{E} [\max (1 - D_\psi(\mathbf{X}_i^s, 0))] + \mathbb{E} [\max (1 + D_\psi(G_\theta(\mathbf{X}_i^s, \mathcal{A})), 0)]. \quad (2)$$

While D_ψ promotes real-looking images, it does not preserve the content or the calligraphic styles. To preserve the textual content in the generated samples we use a handwritten recognizer network R_ϕ that examines whether the gen-

erated samples are actually real text. The recognizer R_ϕ is inspired by CRNN [21]. The CTC loss [10] is used to compare the recognizer output to the query words that were given as input to G_θ . Recognizer R_ϕ is only optimized with real, labelled, handwritten samples, but it is used to encourage G_θ to produce readable text with accurate content. The loss is defined as,

$$L_R = \mathbb{E}_{\mathbf{x} \sim \{\mathbf{X}_i^s, \tilde{\mathbf{X}}_i^t\}} \left[- \sum \log(p(y_r | R_\phi(\mathbf{x}))) \right]. \quad (3)$$

Here, y_r is the transcription string of $\mathbf{x} \sim \{\mathbf{X}_i^s, \tilde{\mathbf{X}}_i^t\}$.

A style classifier network S_η is employed to guide the network G_θ in producing samples conditioned to a particular writing style. The network S_η attempts to predict the writer of a given handwritten image. The cross-entropy objective is applied as a loss function. S_η is trained only on the real samples using the loss given below,

$$L_S = \mathbb{E}_{\mathbf{x} \sim \{\mathbf{X}_i^s, \tilde{\mathbf{X}}_i^t\}} \left[- \sum y_i \log(S_\eta(\mathbf{x})) \right]. \quad (4)$$

An important feature of our design is to utilize a cycle loss that ensures the encoded style features have cycle consistency. This loss function enforces the decoder to preserve the style information in the decoding process, such that the original style feature sequence can be reconstructed from the generated image. Given the generated word images $\tilde{\mathbf{X}}_i^t$, we use the encoder $T_\mathcal{E}$ to reconstruct the style feature sequence $\tilde{\mathbf{Z}}$. The cycle loss L_c minimizes the error between the style feature sequence \mathbf{Z} and its reconstruction $\tilde{\mathbf{Z}}$ by means of a L_1 distance metric,

$$L_c = \mathbb{E} \left[\left\| T_\mathcal{E}(\mathbf{X}_i^s) - T_\mathcal{E}(\tilde{\mathbf{X}}_i^t) \right\|_1 \right]. \quad (5)$$

The cycle loss imposes a regularization to the decoder for consistently imitating the writing style in the generated styled text images. Overall, we train our HWT model in an end-to-end manner with the following loss objective,

$$L_{total} = L_{adv} + L_S + L_R + L_c. \quad (6)$$

We observe balancing the gradients of the network S_η and R_ϕ is helpful in the training with our loss formulation. Following [3], we normalize the ∇S_η and ∇R_ϕ to have the same standard deviation (σ) as adversarial loss gradients,

$$\nabla S_\eta \leftarrow \alpha \left(\frac{\sigma_D}{\sigma_S} \cdot \nabla S_\eta \right), \nabla R_\phi \leftarrow \alpha \left(\frac{\sigma_D}{\sigma_R} \cdot \nabla R_\phi \right). \quad (7)$$

Here, α is a hyper-parameter that is fixed to 1 during the training of our model.

4. Experiments

We perform extensive experiments on IAM handwriting dataset [19]. It consists of 9862 text lines with around

Table 1: **Comparison of the HWT with GANwriting [14] and Davis *et al.* [5]** in terms of FID scores computed between the generated text images and real text images of the IAM dataset. Our HWT performs favorably against [14, 5] in all four settings: In-Vocabulary words and seen style (IV-S), In-Vocabulary words and unseen style (IV-U), Out-of-vocabulary content and seen style (OOV-S) and Out-of-vocabulary content and unseen style (OOV-U). On the challenging setting of OOV-U, HWT achieves an absolute gain of 16.5 in FID score, compared to GANwriting [14]. Best results are in bold.

	IV-S ↓	IV-U ↓	OOV-S ↓	OOV-U ↓
GANwriting [14]	120.07	124.30	125.87	130.68
Davis <i>et al.</i> [5]	118.56	128.75	127.11	136.67
HWT (Ours)	106.97	108.84	109.45	114.10

62,857 English words, written by 500 different writers. For thorough evaluation, we reserve an exclusive subset of 160 writers for testing, while images from the remaining 340 writers are used for our model training. In all our experiments, we resize images to a fixed height of 64 pixels, while maintaining the aspect ratio of original image. For training, we use $P = 15$ style example images, as in [14]. Both the transformer encoder and transformer decoder networks employ three attention layers ($L = 3$) and each attention layer applies multi-headed attention having 8 attention heads ($J = 8$). We set the embedding size d to 512. In all experiments, we train our model for 4k epochs with a batch size of 8 on a single V100 GPU. Adam optimizer is employed during training with a learning rate of 0.0002.

4.1. Styled Handwritten Text Generation

We first evaluate (Tab. 1) our approach for styled handwritten text image generation, where both style and content are desired to be imitated in the generated text image. Following [14], we use Fr chet Inception Distance (FID) [13] evaluation metric for comparison. The FID metric is measured by computing the distance between the Inception-v3 features extracted from generated and real samples for each writer and then averaging across all writers. We evaluate our HWT with GANwriting [14] and Davis *et al.* [5] in four different settings: In-Vocabulary words and seen styles (IV-S), In-Vocabulary words and unseen styles (IV-U), Out-of-Vocabulary words and seen styles (OOV-S), and Out-of-Vocabulary words and unseen styles (OOV-U). Among these settings, most challenging one is the OOV-U, where both words and writing styles are never seen during training. For OOV-S and OOV-U settings, we use a set of 400 words that are distinct from IAM dataset transcription, as in [14]. In all four settings, the transcriptions of real samples and generated samples are different. Tab. 1 shows that HWT performs favorably against both existing methods [14, 5].

Fig 4 presents the qualitative comparison of HWT with

Style examples	HWT (Ours)	GANwriting	Davis et al.
<i>A good neighbour to those Africans who will not think to live in houses of wood and stone of</i> The process has been too slow for Herr Strauss and last month he attacked Britain for being an There were loud cries of 'shame' from all parts of the Conservative side Mr. Hill appeared to be in He thought he said, 7 of the Soviet Union would be prepared to reach an agreement for the zone of Mr. Macleod went on with the conference at Lancaster House despite the crisis which had blown By the end of the month he still delighted in Naples he told Clamancy that he enjoyed it all	<i>No two people can write precisely the same way just like no two people can have the same fingerprints</i> No two people can write precisely the same way just like no two people can have the same fingerprints No two people can write precisely the same way just like no two people can have the same fingerprints No two people can write precisely the same way just like no two people can have the same fingerprints No two people can write precisely the same way just like no two people can have the same fingerprints No two people can write precisely the same way just like no two people can have the same fingerprints No two people can write precisely the same way just like no two people can have the same fingerprints	<i>No two people can write precise the same way just like no two people can have the same fingerprints</i> No two people can write precise the same way just like no two people can have the same fingerprints No two people can write precise the same way just like no two people can have the same fingerprints No two people can write precise the same way just like no two people can have the same fingerprints No two people can write precise the same way just like no two people can have the same fingerprints No two people can write precise the same way just like no two people can have the same fingerprints No two people can write precise the same way just like no two people can have the same fingerprints	<i>No two people can write precisely the same way, just like no two people can have the same fingerprints</i> No two people can write precisely the same way, just like no two people can have the same fingerprints No two people can write precisely the same way, just like no two people can have the same fingerprints No two people can write precisely the same way, just like no two people can have the same fingerprints No two people can write precisely the same way, just like no two people can have the same fingerprints No two people can write precisely the same way, just like no two people can have the same fingerprints No two people can write precisely the same way, just like no two people can have the same fingerprints

Figure 4: Qualitative comparison of our HWT (second column) with GANwriting [14] (third column) and Davis et al. [5] (fourth column). We use the same textual content 'No two people can write precisely the same way just like no two people can have the same fingerprints' for all three methods. The first column shows the style examples from different writers. Davis et al. [5] captures the global style, e.g. slant, but struggles to mimic the character-specific style details. On the other hand, since GANwriting [14] is limited to a fixed length query words, it is unable to complete the provided textual content. Our HWT better mimics global and local style patterns, generating more realistic handwritten text images.

[14, 5] for styled handwritten text generation. We present results for different writers, whose example style images are shown in the first column. For all the three methods, we use the same textual content. While Davis et al. [5] follows the leftward slant of the last style example from the top, their approach struggles to capture character-level styles and curvilinear patterns (e.g. see the word 'the'). On the other hand, GANwriting [14] struggles to follow leftward slant of the last style example from the top and character-level styles. Our HWT better imitates both the global and local style patterns in these generated example text images.

4.2. Handwritten Text Generation

Here, we evaluate the quality of the handwritten text image generated by our HWT. For a fair comparison with the recently introduced ScrabbleGAN [7] and Davis et al. [5], we report our results in the same evaluation settings as used by [7, 5]. Tab. 2 presents the comparison with [7, 5] in terms of FID and geometric-score (GS). Our HWT achieves favourable performance, compared to both approaches in terms of both FID and GS scores. Different from Tab. 1, the results reported here in Tab. 2 indicates the quality of the generated images, compared with the real examples in the IAM dataset, while ignoring style imitation capabilities.

4.3. Ablation study

We perform multiple ablation studies on the IAM dataset to validate the impact of different components in our frame-

Table 2: **Handwritten text image generation quality comparison of our proposed HWT with ScrabbleGAN [7] and Davis et al. [5] on the IAM dataset.** Results are reported in terms of FID and GS by following the same evaluation settings, as in [7, 5]. Our HWT performs favorably against these methods in terms of both FID and GS. Best results are in bold.

	FID ↓	GS ↓
ScrabbleGAN [7]	20.72	2.56×10^{-2}
Davis et al. [5]	20.65	4.88×10^{-2}
HWT (Ours)	19.40	1.01×10^{-2}

work. Tab. 3 shows the impact of integrating transformer encoder (Enc), transformer decoder (Dec) and cycle loss (CL) to the baseline (Base). Our baseline neither uses transformer modules nor utilizes cycle loss. It only employs a CNN encoder to obtain style features, whereas the content features are extracted from the one-hot representation of query words. Both content and style features are passed through a CNN decoder to generate styled handwritten text images. While the baseline is able to generate realistic text images, it has a limited ability to mimic the given writer's style leading to inferior FID score (row 1). The introduction of the transformer encoder into the baseline (row 2) leads to an absolute gain of 5.6 in terms of FID score, highlighting the importance of our transformer-based self-attentive feature sequence in the generator encoder. We ob-

Table 3: **Impact of integrating transformer encoder (Enc), transformer decoder (Dec) and cycle loss (CL) to the baseline (Base)** on the OOV-U settings of IAM dataset. Results are reported in terms of FID score. Best results are reported in bold. On right, we show the effect of each component when generating two example words ‘freedom’ and ‘precise’ mimicking two given writing styles.

	FID ↓	Style Example
		<i>admired people</i>
Base	134.45	<i>freedom precise</i>
Base + Enc	128.80	<i>freedom precise</i>
Base + Dec	124.81	<i>freedom precise</i>
Base + Enc + Dec	116.50	<i>freedom precise</i>
Base + Enc + Dec + CL	114.10	<i>freedom precise</i>

serve here that the generated sample still lacks details in terms of character-specific style patterns. When integrating the transformer decoder into the baseline (row 3), we observe a significant gain of 9.6 in terms of FID score. Notably, we observe a significant improvement (17.9 in FID) when integrating both transformer encoder and decoder to the baseline (row 4). This indicates the importance of self-and encoder-decoder attention for achieving realistic styled handwritten text image generation. The performance is further improved by the introduction of cycle loss to our final HWT architecture (row 4).

As described earlier (Sec. 3.2), HWT strives for style-content entanglement at character-level by feeding query character embeddings to the transformer decoder network. Here, we evaluate the effect of character-level content encoding (conditioning) by replacing it with word-level conditioning. We obtain the word-level embeddings, by using an MLP that aims to obtain string representation of each query word. These embeddings are used as conditional input to the transformer decoder. Table 4 suggests that HWT benefits from character-level conditioning that ensures finer control of text style. The performance of word-level conditioning is limited to mimicking the global style, whereas our character-level approach ensures locally realistic as well as globally consistent style patterns.

4.4. Human Evaluation

Here, we present results of our two user studies on 100 human participants to evaluate whether the proposed HWT achieves human plausibility in terms of the style mimicry. First, a *User preference study* compares styled text images generated by our method with GANwriting [14] and Davis *et al.* [5]. Second, a *User plausibility study* that evaluates the proximity of the synthesized samples generated by our method to the real samples. In both studies, synthesized samples are generated using *unseen writing styles* of test set

Table 4: **Comparison between word and character-level conditioning** on IAM dataset. Results are reported in terms of FID score. Our character-level conditioning performs favorably, compared to its word-level counterpart. Best results are reported in bold. On the right, we show the effect of word and character-level conditioning, when generating two example words ‘symbols’ and ‘same’ mimicking two given writing styles.

	FID ↓	Style Example
		<i>Liberty have</i>
Word-level	126.87	<i>symbols same</i>
Character-level	114.10	<i>symbols same</i>

writers of IAM dataset, and for textual content we use sentences from Stanford Sentiment Treebank [22] dataset.

For *User preference study*, each participant is shown the real handwritten paragraph of a person and synthesized handwriting samples of that person using HWT, Davis *et al.* [5] and GANwriting [14], randomly organized. The participants were asked to mark the best method for mimicking the real handwriting style. In total, we have collected 1000 responses. The results of this study shows that our proposed HWT was preferred 81% of the time over the other two methods.

For *User plausibility study*, each participant is shown a person’s actual handwriting, followed by six samples, where each of these samples is either genuine or synthesized handwriting of the same person. Participants are asked to identify whether a given handwritten sample is genuine or not (forged/synthesized) by looking at the examples of the person’s real handwriting. Thus, each participant provides 60 responses, thereby we collect 6000 responses for 100 participants. For this study, only 48.1% of the images have been correctly classified, thereby showing a comparable performance to a random choice in a two-class problem.

5. Conclusion

We introduced a transformer-based styled handwritten text image generation approach, HWT, that comprises a conditional generator having an encoder-decoder network. Our HWT captures the long and short range contextual relationships within the writing style example through a self-attention mechanism, thereby encoding both global and local writing style patterns. In addition, HWT utilizes an encoder-decoder attention that enables style-content entanglement at the character-level by inferring the style representation for each query character. Qualitative, quantitative and human-based evaluations show that our HWT produces realistic styled handwritten text images with varying length and any desired writing style.

References

- [1] Emre Aksan and Otmar Hilliges. Stcn: Stochastic temporal convolutional networks. *arXiv preprint arXiv:1902.06568*, 2019. 2
- [2] Emre Aksan, Fabrizio Pece, and Otmar Hilliges. Deepwriting: Making digital ink editable via deep generative modeling. In *CHI*, pages 1–14, 2018. 2
- [3] Eloi Alonso, Bastien Moysset, and Ronaldo Messina. Adversarial generation of handwritten text images conditioned on sequences. In *ICDAR*, pages 481–486. IEEE, 2019. 1, 2, 3, 6
- [4] Bo Chang, Qiong Zhang, Shenyi Pan, and Lili Meng. Generating handwritten chinese characters using cyclegan. In *WACV*, pages 199–207. IEEE, 2018. 1, 2, 3
- [5] Brian Davis, Chris Tensmeyer, Brian Price, Curtis Wigington, Bryan Morse, and Rajiv Jain. Text and style conditioned gan for generation of offline handwriting lines. *BMVC*, 2020. 1, 2, 3, 6, 7, 8, 5, 9, 10, 11, 12, 13
- [6] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*, 2020. 5
- [7] Sharon Fogel, Hadar Averbuch-Elor, Sarel Cohen, Shai Mazor, and Roei Litman. Scrabblegan: semi-supervised varying length handwritten text generation. In *CVPR*, pages 4324–4333, 2020. 1, 2, 3, 7
- [8] Ian J Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial networks. *arXiv preprint arXiv:1406.2661*, 2014. 1
- [9] Alex Graves. Generating sequences with recurrent neural networks. *arXiv preprint arXiv:1308.0850*, 2013. 2
- [10] Alex Graves, Santiago Fernández, Faustino Gomez, and Jürgen Schmidhuber. Connectionist temporal classification: labelling unsegmented sequence data with recurrent neural networks. In *ICML*, pages 369–376, 2006. 6
- [11] Tom SF Haines, Oisín Mac Aodha, and Gabriel J Brostow. My text in your handwriting. *TOG*, 35(3):1–18, 2016. 3
- [12] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *CVPR*, pages 770–778, 2016. 5
- [13] Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. Gans trained by a two time-scale update rule converge to a local nash equilibrium. *arXiv preprint arXiv:1706.08500*, 2017. 6
- [14] Lei Kang, Pau Riba, Yaxing Wang, Marçal Rusiñol, Alicia Fornés, and Mauricio Villegas. Ganwriting: Content-conditioned generation of styled handwritten word images. In *ECCV*, pages 273–289. Springer, 2020. 1, 2, 3, 6, 7, 8, 5, 9, 10, 11, 12
- [15] Tero Karras, Samuli Laine, and Timo Aila. A style-based generator architecture for generative adversarial networks. In *CVPR*, pages 4401–4410, 2019. 2
- [16] Salman Khan, Muzammal Naseer, Munawar Hayat, Syed Waqas Zamir, Fahad Shabbaz Khan, and Mubarak Shah. Transformers in vision: A survey. *arXiv preprint arXiv:2101.01169*, 2021. 4
- [17] Atsunobu Kotani, Stefanie Tellex, and James Tompkin. Generating handwriting via decoupled style descriptors. In *ECCV*, pages 764–780. Springer, 2020. 2
- [18] Jae Hyun Lim and Jong Chul Ye. Geometric gan. *arXiv preprint arXiv:1705.02894*, 2017. 5
- [19] U-V Marti and Horst Bunke. The iam-database: an english sentence database for offline handwriting recognition. *IJ-DAR*, 5(1):39–46, 2002. 6
- [20] Martin Mayr, Martin Stumpf, Angelos Nikolaou, Mathias Seuret, Andreas Maier, and Vincent Christlein. Spatio-temporal handwriting imitation. *arXiv preprint arXiv:2003.10593*, 2020. 3
- [21] Baoguang Shi, Xiang Bai, and Cong Yao. An end-to-end trainable neural network for image-based sequence recognition and its application to scene text recognition. *PAMI*, 39(11):2298–2304, 2016. 6
- [22] Richard Socher, Alex Perelygin, Jean Wu, Jason Chuang, Christopher D Manning, Andrew Y Ng, and Christopher Potts. Recursive deep models for semantic compositionality over a sentiment treebank. In *EMNLP*, pages 1631–1642, 2013. 8, 1
- [23] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, undefinedukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *NIPS*, page 6000–6010, Red Hook, NY, USA, 2017. Curran Associates Inc. 5
- [24] Jun-Yan Zhu, Taesung Park, Phillip Isola, and Alexei A Efros. Unpaired image-to-image translation using cycle-consistent adversarial networks. In *ICCV*, pages 2223–2232, 2017. 3

Handwriting Transformers

Supplementary Material

In this supplementary material, we present additional human study details, additional qualitative results, and additional ablation study results. In Sec. 1, we provide details of human study experiments. Sec. 2 presents the additional visualisation results of transformer encoder-decoder attention maps. Sec. 3 shows qualitative comparison of our proposed HWT. Sec. 4 shows the interpolations between two different calligraphic styles on the IAM dataset. Finally, Sec. 5 presents additional ablation results.

1. Human Study Additional Details

Here, we present results of our two user studies on 100 human participants to evaluate the human plausibility in terms of style mimicry of our proposed HWT. In both these user studies, the forged samples are generated using *unseen writing styles* of test set writers of IAM dataset, and for textual content we use sentences from Stanford Sentiment Treebank [22] dataset.

User Preference Study: Fig. 1 shows the interface for the *User preference study* experiment, which compares styled text images. In this study, each participant is shown a real handwritten text image of a person and the synthesized handwriting text images of that person using our proposed HWT, Davis *et al.* [5] and GANwriting [14]. We randomly present generated results of these methods to the user. Then, the user can compare the real image and the generated images side by side on the same screen and without any time restriction to give the answer. Each participant is required to provide response for a total of ten questions. Overall, we have collected 1000 responses from 100 participants. Table 1 shows the results of *User preference study*. Davis *et al.* [5] and GANwriting [14] were preferred 9% (90 responses out of the total 1000) and 10% (100 responses out of the total 1000), respectively. Our proposed HWT was preferred 81% (810 responses out of the total 1000 responses) over the other two existing methods.

User Plausibility Study: Fig. 2 shows the interface for the *User plausibility study*, which evaluates the proximity of the synthesized samples generated by our proposed HWT to the real samples. Here, each participant is shown a person’s actual handwriting, followed by six samples, where each of these samples is either genuine or synthesized handwriting of the same person. Participants are asked to identify whether a given handwritten sample is genuine or not (forged/synthesized) with no time limit restriction to answer the question. In total, we collect 6000 responses for 100 human participants as each one provides 60 responses.

Table 1: *User preference study* in comparison to GANwriting [14] and Davis *et al.* [5]. The result shows that our proposed HWT was preferred 81% of the time over the other two compared methods.

	Total Responses	User Preferences
GANwriting [14]	1000	100
Davis <i>et al.</i> [5]		90
HWT (Ours)		810

Table 2: Confusion matrix (%) obtained from *User plausibility study*. Only 48.1% of the images were correctly classified, indicating an output comparable to a random choice in a two-class problem.

Actual	Predicted	
	Real	Fake
Real	24.9	25.1
Fake	26.8	23.2

The study revealed that the generated images produced by our proposed HWT were deemed plausible. Table 2 shows the confusion matrix of the human assessments. For this particular study, only 48.1% of the images have been correctly classified, which indicates a comparable performance to random choice in a two-class problem.

2. Additional Visualizations of Transformer Encoder-Decoder Attention

Fig. 3 shows the visualization of attention maps obtained using encoder-decoder of our approach (HWT) at the last layer of the transformer decoder. We compute the attention matrices for four different words: ‘laughs’, ‘because’, ‘inside’, and ‘fashion’. Note that the attention maps generated by our model focus on the relevant regions of interest in the style examples for each query character. For instance, to infer character-specific style attributes of a given character ‘h’ in the query word ‘laughs’, the model gives priority to multiple image regions containing the character ‘h’. Note that if the query character isn’t found in the style examples, the model attempts to find similar characters. For example, to obtain character representation of ‘u’ in the query word ‘laughs’, the attention algorithm highlights image regions containing similar characters (*e.g.* ‘n’).

Figure 1: Screenshot of the Interface used in *User preference study* experiment. Each participant is shown the real handwritten text image (on the left side) of a person and synthesized handwriting text images (on the right side) of that person generated using three different methods. Participants have to mark the best method for mimicking the real handwriting style.

User preference study: Instructions

1. The image on the left side shows a real example of handwriting style of a person. 2. The images on the right side are generated by computer using three different methods (randomly ordered). Please note that the textual content in the right images is different from the left image. Question: Which one is better at mimicking the handwriting style of the left image? Please give your response by clicking on the checkboxes.

Geac searched in his pockets once more and came
paper-clips After a few seconds of hoisting
a bent wire loop into the lock and rattled it

- ☐ Shaky closeups of turkeyonrolls stubbly chins liver
spots red noses and the filmmakers new bobbed do
draw easy chuckles but lead
- ☐ Shaky closeups of turkeyonrolls stubbly
chins liver spots red noses and the filmmakers
new bobbed do draw easy chuckles but lead
- ☐ Shaky closeup of turkeyo stubbly chins
spots red nose and the filmmak new bobbed do
draw easy chuckle but lead

Next [0/10]

3. Additional Qualitative Comparison

Figs. 4-21 show qualitative comparison between our proposed HWT with [14, 5] for styled handwritten text generation. Note that we use the same textual content for all the examples figures for all the three methods to ensure a fair comparison. The first row in each figure presents the different writers example style images. The rest of the rows correspond to our HWT and [14, 5] respectively. The qualitative results suggest that our method is promising at imitating character-level patterns, while the other two methods struggle to retain character-specific details. The success of the other two methods is limited to capturing only the global patterns (e.g., slant, ink widths). In some cases, these methods even struggle to capture global styles. In Fig. 6, Fig. 16 and Fig. 18, Davis *et al.* [5] suffer to capture the slant. Whereas, in Fig. 16 and Fig. 20, the ink width of the images generated by this method is not consistent with the style examples. On the other hand, since GANwriting [14] is limited to a fixed length query words, it is unable to complete few words that exceed the limit.

Figs. 22-23 show qualitative results using the same text as in the style examples to compare our proposed HWT with [14, 5]. Figs. 24-26 show examples, where we aim to gener-

ate arbitrarily long words. The results show that our model is capable of consistently imitating the styles of the given style example, even for arbitrarily long words. Note that GANwriting [14] struggles to generate long words.

4. Latent Space Interpolations

Fig. 27 shows interpolations between two different calligraphic styles on the IAM dataset. To interpolate by λ between two writers A and B , we compute the weighted average $Z_{AB} = \lambda Z_A + (1 - \lambda) Z_B$, while keeping the textual contents fixed. Here, Z_A and Z_B are the style feature sequences obtained from encoder T_E . It is worth mentioning that our models produce images seamlessly by adjusting from one style to other, which indicates that our model generalizes in the latent space rather than memorizing any trivial writing patterns.

5. Additional Ablation Results

Fig. 28 presents additional qualitative results that show the impact of integrating transformer encoder (Enc), transformer decoder (Dec) and cycle loss (CL) to the baseline (Base). Fig. 29 shows additional qualitative comparisons between word-level and character-level conditioning.

Figure 2: Screenshot of the Interface used in *User plausibility study* experiment. Each participant is shown a person's actual handwriting (on the left side), followed by six samples (on the right side), where three out of these samples are genuine and the rest are synthesized. Participants have to classify each sample as genuine or forgery by looking at the real image.

User plausibility study: Instructions

1. The image on the left side shows a real example of handwriting style of a person. 2. The images on the right side are either real handwriting or forged handwriting of that person. In total there are three real images and three forged images. You have to identify the forged samples out of these images.

that sinister yarn in which half the cast try to persuade

☐ heroine that she is out of her mind Despite flagrant cheating

☐ Director Andrew Niccol ... demonstrates a very understanding of the quirks of fame

☐ Happy Times maintains an appealing veneer without becoming too cute about it

☐ He has earned his break The film is a well-made variation

☐ eerie atmosphere is built up neatly Susan Strassberg is the crippled

☐ A yarn that respects the Marvel version without becoming ensnared by it

Next [0/10]

Figure 3: Additional visualization results of encoder-decoder attention maps at the last layer of the transformer decoder. The attention maps are computed for four different query words: ‘laughs’, ‘because’, ‘inside’, and ‘fashion’. Heat maps corresponding to all characters (including repetitions, as the letter ‘i’ appears twice in ‘inside’) of these query words are shown in the figure.



Figure 4: Additional qualitative comparisons of our proposed HWT with GANwriting [14] and Davis et al. [5], when generating the same text ‘With more character development this might have been an eerie thriller with better payoffs it could have been a thinking’.

Style	<i>sorts of lessons or you might. 'leave him' But how utterly societies for that. 'mother will allow for the wife's</i>
HWT (Ours)	<i>With more character development this might have been an eerie thriller with better payoffs it could have been a thinking</i>
GANwriting	<i>With more charact develop this might have been an eerie thriller with better payoffs it could have been a thinkin</i>
Davis et al.	<i>With more character development this might have been an eerie thriller with better payoffs it could have been a thinking</i>

Figure 5: Additional qualitative comparisons of our proposed HWT with GANwriting [14] and Davis et al. [5], when generating the same text ‘Its not helpful to listen to extremist namecalling regardless of whether you think Kissinger was a calculating’.

Style	<i>between the do-it-yourself cupboard and the polished brass fourteen-pounder-shell-case which served respectively as cool</i>
HWT (Ours)	<i>It's not helpful to listen to extremist namecalling regardless of whether you think Kissinger was a calculating</i>
GANwriting	<i>It's not helpful to listen to extremi namecal regardl of whether you think kissing was a calcula</i>
Davis et al.	<i>It's not helpful to listen to extremist namecalling regardless of whether you think Kissinger was a calculating</i>

Figure 6: Additional qualitative comparisons of our proposed HWT with GANwriting [14] and Davis et al. [5], when generating the same text ‘Shaky closeups of turkeyonrolls stubbly chins liver spots red noses and the filmmakers new bobbed do draw easy chuckles but’.

Style	<i>Ceac searched in his pockets once more and came up with two paper-clips after a few seconds of hoisting the roughly thrust</i>
HWT (Ours)	<i>Shaky closeups of turkeyonrolls stubbly chins liver spots red noses and the filmmakers new bobbed do draw easy chuckles but</i>
GANwriting	<i>Shaky closeup of turkeyo stubbly chins liver spots red noses and the filmmak new bobbed do draw easy chuckle but</i>
Davis et al.	<i>Shaky closeups of turkeyonrolls stubbly chins liver spots red noses and the filmmakers new bobbed do draw easy chuckles but</i>

Figure 7: Additional qualitative comparisons of our proposed HWT with GANwriting [14] and Davis *et al.* [5], when generating the same text ‘This film was made by and for those folks who collect the serial killer cards and are fascinated by the mere suggestion’.

Style	He read the film star's sorry story and frowned at the provisions of Schedule D taxation which not only allowed his
HWT (Ours)	This film was made by and for those folks who collect the serial killer cards and are fascinated by the mere suggestion
GANwriting	This film was made by and for those folks who collect the serial killer cards and are fascina by the mere suggest
Davis <i>et al.</i>	This film was made by and for those folks who collect the serial killer cards and are fascinated by the mere suggestion

Figure 8: Additional qualitative comparisons of our proposed HWT with GANwriting [14] and Davis *et al.* [5], when generating the same text ‘Its a drawling slobbering lovable runon sentence of a film a Southern Gothic with the emotional arc of its raw blues’.

Style	Compared with 1958 the expenditure index for 1959 showed increases of 4 to 6 per cent for couples without children and
HWT (Ours)	It is a drawling slobbering lovable runon sentence of a film Southern Gothic with the emotional arc of its raw blues
GANwriting	It is a drawlin slobber lovable runon sentence of a film Southern Gothic with the emotion arc of its raw blues
Davis <i>et al.</i>	It is a drawling slobbering lovable runon sentence of a film a Southern Gothic with the emotional arc of its raw blues

Figure 9: Additional qualitative comparisons of our proposed HWT with GANwriting [14] and Davis *et al.* [5], when generating the same text ‘LRB W RRB hile long on amiable monkeys and worthy environmentalism Jane Goodalls Wild Chimpanzees is short’.

Style	In April of that year his first wife's brother-in-law the diplomatist Lord Ponsonby had written to advise Annesley
HWT (Ours)	LRB W RRB hile long on amiable monkeys and worthy environmentalism Jane Goodall's Wild Chimpanzees is short
GANwriting	LRB W RRB hile long on amiable monkeys and worthy environ Jane Goodall's Wild Chimpan is short
Davis <i>et al.</i>	LRB W RRB hile long on amiable monkeys and worthy environmentalism Jane Goodall's Wild Chimpanzees is short

Figure 10: Additional qualitative comparisons of our proposed HWT with GANwriting [14] and Davis *et al.* [5], when generating the same text ‘For close to two hours the audience is forced to endure three terminally depressed mostly inarticulate hyper dysfunctional’

Style	give the system a trial adding that it was being cultivated with extraordinary success in France and Italy and that he
HWT (Ours)	For close to two hours the audience is forced to endure three terminally depressed mostly inarticulate hyper dysfunctional
GANwriting	For close to two hours the audience is forced to endure three terminally depressed mostly inarticulate hyper dysfunctional
Davis <i>et al.</i>	For close to two hours the audience is forced to endure three terminally depressed mostly inarticulate hyper dysfunctional

Figure 11: Additional qualitative comparisons of our proposed HWT with GANwriting [14] and Davis *et al.* [5], when generating the same text ‘Claude Chabrol’s camera has a way of gently swaying back and forth as it cradles its characters veiling tension beneath’.

Style	Are there differences in adjustment to ageing and retirement according to the occupational level of employees
HWT (Ours)	Claude Chabrol’s camera has a way of gently swaying back and forth as it cradles its characters veiling tension beneath
GANwriting	Claude Chabrol’s camera has a way of gently swaying back and forth as it cradles its characters veiling tension beneath
Davis <i>et al.</i>	Claude Chabrol’s camera has a way of gently swaying back and forth as it cradles its characters veiling tension beneath

Figure 12: Additional qualitative comparisons of our proposed HWT with GANwriting [14] and Davis *et al.* [5], when generating the same text ‘Though the plot is predictable the movie never feels formulaic because the attention is on the nuances of the’.

Style	None of the numerous conventional remedies to which he had been subjected until since the symptoms had first shown
HWT (Ours)	Though the plot is predictable the movie never feels formulaic because the attention is on the nuances of the
GANwriting	Though the plot is predictable the movie never feels formulaic because the attention is on the nuances of the
Davis <i>et al.</i>	Though the plot is predictable the movie never feels formulaic because the attention is on the nuances of the

Figure 13: Additional qualitative comparisons of our proposed HWT with GANwriting [14] and Davis *et al.* [5], when generating the same text ‘A comingofage tale from New Zealand whose boozy languid air is balanced by a rich visual clarity and deeply felt’.

Style	As a result the Glasgow Retirement Council came into being in April 1958 with Dr. Andrew Hood as chairman and Mr. Andrew
HWT (Ours)	A comingofage tale from New Zealand whose boozy languid air is balanced by a rich visual clarity and deeply felt
GANwriting	A comingo tale from New Zealand whose boozy languid air is balance by a rich visual clarity and deeply felt
Davis et al.	A comingofage tale from New Zealand whose boozy languid air is balanced by a rich visual clarity and deeply felt

Figure 14: Additional qualitative comparisons of our proposed HWT with GANwriting [14] and Davis *et al.* [5], when generating the same text ‘Unfortunately Kapur modernizes AEW. Masons story to suit the sensibilities of a young American a decision that plucks The’.

Style	the problems of ad meaning or already in retirement but also strong desire the part of all concerned, for concerted action
HWT (Ours)	Unfortunately Kapur modernizes AEW Mason's story to suit the sensibilities of a young American a decision that plucks The
GANwriting	Unfortun Kapur moderni AEW Mason's story to suit the sensibi of a young America a decisio that plucks The
Davis et al.	Unfortunately Kapur modernizes AEW Mason's story to suit the sensibilities of a young American a decision that plucks The

Figure 15: Additional qualitative comparisons of our proposed HWT with GANwriting [14] and Davis *et al.* [5], when generating the same text ‘Unless Bob Crane is someone of particular interest to you this films impressive performances and adept direction are’.

Style	The Putney Owen scheme is now in its fourth year and opportunity has been taken to revise the course in the light
HWT (Ours)	Unless Bob Crane is someone of particular interest to you this film's impressive performances and adept direction are
GANwriting	Unless Bob Crane is someone of particu interes to you this film's impress perform and adept directi are
Davis et al.	Unless Bob Crane is someone of particular interest to you this film's impressive performances and adept direction are

Figure 16: Additional qualitative comparisons of our proposed HWT with GANwriting [14] and Davis *et al.* [5], when generating the same text ‘Affirms the gifts of all involved starting with Spielberg and going right through the ranks of the players oncamera and off’.

Style	nevertheless identify it. Similarly the psychologist has to be prepared to observe and make inferences about all kinds of
HWT (Ours)	Affirms the gifts of all involved starting with Spielberg and going right through the ranks of the players oncamera and off
GANwriting	Affirms the gifts of all involve startin with Spielbe and going right through the ranks of the players oncamer and off
Davis <i>et al.</i>	Affirms the gifts of all involved starting with Spielberg and going right through the ranks of the players oncamera and off

Figure 17: Additional qualitative comparisons of our proposed HWT with GANwriting [14] and Davis *et al.* [5], when generating the same text ‘Though this rude and crude film does deliver a few gut-busting laughs its digs at modern society are all things we ve seen’.

Style	headlines for having smashed large number & to relieve his feelings On the small-to-medium establishment it is
HWT (Ours)	Though this rude and crude film does deliver a few gutbusting laughs its digs at modern society are all things we ve seen
GANwriting	Though this rude and crude film does deliver a few gutbust laughs its digs at modern society are all things we ve seen
Davis <i>et al.</i>	Though this rude and crude film does deliver a few gutbusting laughs its digs at modern society are all things we ve seen

Figure 18: Additional qualitative comparisons of our proposed HWT with GANwriting [14] and Davis *et al.* [5], when generating the same text ‘You ll laugh at either the obviousness of it all or its stupidity or maybe even its inventiveness but the point is’.

Style	He had long sensed injustice in the distinctions drawn between ordinary wage-earners and those self-employed By
HWT (Ours)	You ll laugh at either the obviousness of it all or its stupidity or maybe even its inventiveness but the point is
GANwriting	lou ll laugh at either the obvious of it all or its stupidi or maybe even its inventi but the point is
Davis <i>et al.</i>	You ll laugh at either the obviousness of it all or its stupidity or maybe even its inventiveness but the point is

Figure 19: Additional qualitative comparisons of our proposed HWT with GANwriting [14] and Davis *et al.* [5], when generating the same text ‘Writerdirector s Mehta s effort has tons of charm and the whimsy is in the mixture the intoxicating masala of cultures’.

Style	avoided by the installation of a power tool As the element being example one can to enjoy ripping out parts from hardwood that
HWT (Ours)	Writerdirector s Mehta s effort has tons of charm and the whimsy is in the mixture the intoxicating masala of cultures
GANwriting	Writerd s Mehta s effort has tons of charm and the whimsy is in the mixture the intoxic masala of culture
Davis <i>et al.</i>	Writerdirector s Mehta s effort has tons of charm and the whimsy is in the mixture the intoxicating masala of cultures

Figure 20: Additional qualitative comparisons of our proposed HWT with GANwriting [14] and Davis *et al.* [5], when generating the same text ‘While easier to sit through than most of Jaglom s selfconscious and gratingly irritating films it s’.

Style	mysteriously returned to the shed their books tidied and over their shoes cleaned ? How showed their this morning with box
HWT (Ours)	While easier to sit through than most of Jaglom s selfconscious and gratingly irritating films it s
GANwriting	While easier to sit through than most of Jaglom s selfcon and gratingly irritat films it s
Davis <i>et al.</i>	While easier to sit through than most of Jaglom s selfconscious and gratingly irritating films it s

Figure 21: Additional qualitative comparisons of our proposed HWT with GANwriting [14] and Davis *et al.* [5], when generating the same text ‘The connected stories of Breitbart and Hanussen are actually fascinating but the filmmaking in Invincible is such that the’.

Style	to carry through Tory policy Gaitskell's stupid hope The tragedy is that enormous inroads could already have been made
HWT (Ours)	The connected stories of Breitbart and Hanussen are actually fascinating but the filmmaking in Invincible is such that the
GANwriting	The connect stories of Breitba and Hanusse are actual fascina but the filmmak in Invinci is such that the
Davis <i>et al.</i>	The connected stories of Breitbart and Hanussen are actually fascinating but the filmmaking in Invincible is such that the

Figure 22: Reconstruction results using the proposed HWT in comparison to GANwriting [14] and Davis *et al.* [5]. We use the same text as in the style examples to generate handwritten images.

Style	<i>few people in the bar elderly well-off artistic who you felt had made</i>
HWT (Ours)	<i>few people in the bar elderly well-off artistic who you felt had made</i>
GANwriting	<i>few people in the bar elderly well-off artist who you felt had made</i>
Davis <i>et al.</i>	<i>few people in the bar elderly well-off artistic who you felt had made</i>
Style	<i>Fortunately, however the fashion for Victorian architecture which Mr. John Betjeman had started several decades</i>
HWT (Ours)	<i>Fortunately however the fashion for Victorian architecture which Mr. John Betjeman had started several decades</i>
GANwriting	<i>Fortuna however the fashion for Victori archite which Mr John Betjeme had started several decades</i>
Davis <i>et al.</i>	<i>Fortunately however the fashion for Victorian architecture which Mr. John Betjeman had started several decades</i>
Style	<i>families with three child for whom the index declined from index which compares the prices</i>
HWT (Ours)	<i>families with three child for whom the index declined from index which compares the prices</i>
GANwriting	<i>familie with three childre for whom the indec decline from indec which compare the prices</i>
Davis <i>et al.</i>	<i>families with three child for whom the index declined from index which compares the prices</i>
Style	<i>They closed on a single bundle and, fumbling with nervous excitement, he jerked it out.</i>
HWT (Ours)	<i>They closed on a single bundle and fumbling with nervous excitement he pulled it out</i>
GANwriting	<i>They closed on a single bundle and fumblin with nervous excitem he pulled it out</i>
Davis <i>et al.</i>	<i>They closed on a single bundle and fumbling with nervous excitement he pulled it out</i>
Style	<i>Mauro's first action was to write to his revered master at Leipzig asking for</i>
HWT (Ours)	<i>Mauro's first action was to write to his revered master at Leipzig asking for</i>
GANwriting	<i>Mauro's first action was to write to his revered master at Leipzig asking for</i>
Davis <i>et al.</i>	<i>Mauro's first action was to write to his revered master at Leipzig asking for</i>
Style	<i>described his distinguished patient and his symptoms. He told Hahnemann that he</i>
HWT (Ours)	<i>described his distinguished patient and his symptoms He told Hahnemann that he</i>
GANwriting	<i>describ his disting patient and his symptom He told Hahne that he</i>
Davis <i>et al.</i>	<i>described his distinguished patient and his symptoms He told Hahnemann that he</i>
Style	<i>This remarkable of medicine whom Sir Francis Burdett described to Anglesey</i>
HWT (Ours)	<i>This remarkable of medicine whom Sir Francis Burdett described to Anglesey</i>
GANwriting	<i>This remarks of medicin whom Sir Francis Burdett describ to Anglese</i>
Davis <i>et al.</i>	<i>This remarkable of medicine whom Sir Francis Burdett described to Anglesey</i>

Figure 23: Reconstruction results using the proposed HWT in comparison to GANwriting [14] and Davis et al. [5].

Style	Mauro's first action was to write to his revered master at Leipzig asking for
HWT (Ours)	Mauro's first action was to write to his revered master at Leipzig asking for
GANwriting	Mauro's first action was to write to his revered master at Leipzig asking for
Davis et al.	Mauro's first action was to write to his revered master at Leipzig asking for

Style	In April of that year his first wife's brother-in-law the diplomatist
HWT (Ours)	In April of that year his first wife's brother-in-law the diplomatist
GANwriting	In April of that year his first wife's brother the diploma
Davis et al.	In April of that year his first wife's brother-in-law the diplomatist

Style	What a frightful event he wrote I tremble indifferent but I really tremble
HWT (Ours)	What a frightful event he wrote I tremble indifferent but I really tremble
GANwriting	What a frightful event he wrote I tremble indifferent but I really tremble
Davis et al.	What a frightful event he wrote I tremble indifferent but I really tremble

Style	In those early years life was very full both in the parish and in the wider war activities
HWT (Ours)	In those early years life was very full both in the parish and in the wider war activities
GANwriting	In those early years life was very full both in the parish and in the wider war activities
Davis et al.	In those early years life was very full both in the parish and in the wider war activities

Style	He may never have had the disease himself but he can nevertheless identify it
HWT (Ours)	He may never have had the disease himself but he can nevertheless identify it
GANwriting	He may never have had the disease himself but he can nevertheless identify it
Davis et al.	He may never have had the disease himself but he can nevertheless identify it

Style	The large attendance and atmosphere of this Conference held in October
HWT (Ours)	The large attendance and atmosphere of this Conference held in October
GANwriting	The large attendance and atmosphere of this Conference held in October
Davis et al.	The large attendance and atmosphere of this Conference held in October

Style	They are slightly more difficult to manage however until a little
HWT (Ours)	They are slightly more difficult to manage however until a little
GANwriting	They are slightly more difficult to manage however until a little
Davis et al.	They are slightly more difficult to manage however until a little

Figure 24: Handwritten text image generation of arbitrarily long words. We generate the 21-letter word ‘Incomprehensibilities’ in three different styles and compare the results with Davis *et al.* [5].

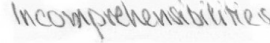
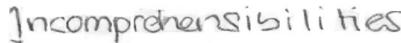


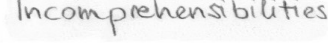
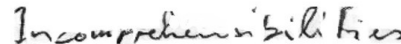
Style	Method	Generated Images
atmosphere	Davis <i>et al.</i>	
	HWT (Ours)	
affiliations	Davis <i>et al.</i>	
	HWT (Ours)	
individuality	Davis <i>et al.</i>	
	HWT (Ours)	

Figure 25: Handwritten text image generation of arbitrarily long words. We generate the 30-letter word ‘Pseudopseudohypoparathyroidism’ in three different styles and compare the results with Davis *et al.* [5].



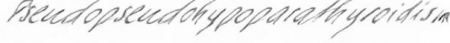

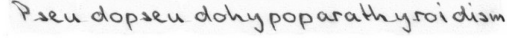
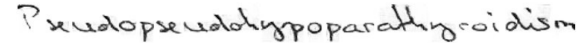
Style	Method	Generated Images
overdraft	Davis <i>et al.</i>	
	HWT (Ours)	
betrayed	Davis <i>et al.</i>	
	HWT (Ours)	
carefully	Davis <i>et al.</i>	
	HWT (Ours)	

Figure 26: Handwritten text image generation of arbitrarily long words. We generate the 28-letter word ‘Antidisestablishmentarianism’ in three different styles and compare the results with Davis *et al.* [5].

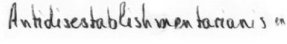
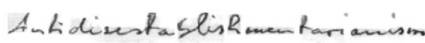
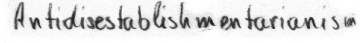

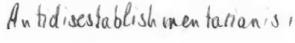
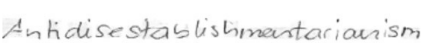
Style	Method	Generated Images
compartment	Davis <i>et al.</i>	
	HWT (Ours)	
delighted	Davis <i>et al.</i>	
	HWT (Ours)	
inspector	Davis <i>et al.</i>	
	HWT (Ours)	

Figure 27: Latent space interpolations between calligraphic styles on the IAM dataset. The first and last image in each column correspond to writing styles of two different writers. Total we have shown five sets of interpolation results. We observe how the generated images seamlessly adjust from one style to another. This result shows that our model can generalize in the latent space rather than memorizing any trivial writing patterns.



Figure 28: Additional qualitative ablation of integrating transformer encoder (Enc), transformer decoder (Dec) and cycle loss (CL) to the baseline (Base) on the IAM dataset. We show the effect of each component when generating six different words ‘especially’, ‘ethereal’, ‘emotional’, ‘standard’, ‘resorts’, and ‘under’.

Style examples →	<i>outwardly</i>	<i>allure</i>	<i>expectations</i>	<i>discovered</i>	<i>reared</i>	<i>content</i>
Base	<i>especially</i>	<i>ethereal</i>	<i>emotional</i>	<i>standard</i>	<i>resorts</i>	<i>under</i>
Base + Enc	<i>especially</i>	<i>ethereal</i>	<i>emotional</i>	<i>standard</i>	<i>resorts</i>	<i>under</i>
Base + Dec	<i>especially</i>	<i>ethereal</i>	<i>emotional</i>	<i>standard</i>	<i>resorts</i>	<i>under</i>
Base + Enc + Dec	<i>especially</i>	<i>ethereal</i>	<i>emotional</i>	<i>standard</i>	<i>resorts</i>	<i>under</i>
Base + Enc + Dec + CL	<i>especially</i>	<i>ethereal</i>	<i>emotional</i>	<i>standard</i>	<i>resorts</i>	<i>under</i>

Figure 29: Additional qualitative comparisons between word and character-level conditioning on IAM dataset. We show the comparison between word and character-level conditioning when generating six different words ‘engaging’, ‘actually’, ‘movie’, ‘rhythms’, ‘what’, and ‘evocative’.

Style examples →	<i>fragile</i>	<i>comically</i>	<i>within</i>	<i>attacked</i>	<i>which</i>	<i>responding</i>
Word-level	<i>engaging</i>	<i>actually</i>	<i>movie</i>	<i>rhythms</i>	<i>what</i>	<i>evocative</i>
Character-level	<i>engaging</i>	<i>actually</i>	<i>movie</i>	<i>rhythms</i>	<i>what</i>	<i>evocative</i>