POLITECNICO DI TORINO Repository ISTITUZIONALE

Self-supervised Text Style Transfer Using Cycle-Consistent Adversarial Networks

Original

Self-supervised Text Style Transfer Using Cycle-Consistent Adversarial Networks / La Quatra, Moreno; Gallipoli, Giuseppe; Cagliero, Luca. - In: ACM TRANSACTIONS ON INTELLIGENT SYSTEMS AND TECHNOLOGY. - ISSN 2157-6912. - ELETTRONICO. - 15:5(2024), pp. 1-38. [10.1145/3678179]

Availability: This version is available at: 11583/2994818 since: 2024-11-28T10:19:00Z

Publisher: ACM

Published DOI:10.1145/3678179

Terms of use:

This article is made available under terms and conditions as specified in the corresponding bibliographic description in the repository

Publisher copyright ACM postprint/Author's Accepted Manuscript

© ACM 2024. This is the author's version of the work. It is posted here for your personal use. Not for redistribution. The definitive Version of Record was published in ACM TRANSACTIONS ON INTELLIGENT SYSTEMS AND TECHNOLOGY, http://dx.doi.org/10.1145/3678179.

(Article begins on next page)

- 4 MORENO LA QUATRA*, Kore University of Enna, Italy
- GIUSEPPE GALLIPOLI*, Politecnico di Torino, Italy
- ⁶ LUCA CAGLIERO, Politecnico di Torino, Italy 7

Text Style Transfer (TST) is a relevant branch of natural language processing that aims to control the style 8 attributes of a piece of text while preserving its original content. To address TST in the absence of parallel 9 data, Cycle-consistent Generative Adversarial Networks (CycleGANs) have recently emerged as promising 10 solutions. Existing CycleGAN-based TST approaches suffer from the following limitations: (1) They apply 11 self-supervision, based on the cycle-consistency principle, in the latent space. This approach turns out to be less 12 robust to mixed-style inputs, i.e., when the source text is partly in the original and partly in the target style; (2) 13 Generators and discriminators rely on recurrent networks, which are exposed to known issues with long-term 14 text dependencies; (3) The target style is weakly enforced, as the discriminator distinguishes real from fake sentences without explicitly accounting for the generated text's style. We propose a new CycleGAN-based TST 15 16 approach that applies self-supervision directly at the sequence level to effectively handle mixed-style inputs and employs Transformers to leverage the attention mechanism for both text encoding and decoding. We also 17 employ a pre-trained style classifier to guide the generation of text in the target style while maintaining the 18 original content's meaning. The experimental results achieved on the formality and sentiment transfer tasks 19 show that our approach outperforms existing ones, both CycleGAN-based and not (including an open-source 20 Large Language Model), on benchmark data and shows better robustness to mixed-style inputs. 21

22 CCS Concepts: • Computing methodologies → Artificial intelligence; • Natural language processing
 23 → Natural language generation;

Additional Key Words and Phrases: Text Style Transfer, Sentiment transfer, Formality transfer, Cycle-consistent
 Generative Adversarial Networks, Transformers

²⁶ ACM Reference Format:

Moreno La Quatra, Giuseppe Gallipoli, and Luca Cagliero. 2024. Self-supervised Text Style Transfer using
 Cycle-Consistent Adversarial Networks. *ACM Trans. Intell. Syst. Technol.* 1, 1 (July 2024), 37 pages. https:
 //doi.org/XXXXXXXXXXXXXXX

1 INTRODUCTION

Language is strongly dependent on both the writers/speakers' characteristics and its context of use (e.g., time, place, scenario, intent). Although humans naturally take these factors into account, Artificial Intelligence systems could struggle to properly handle these aspects. As a result, the development of Natural Language Processing (NLP) tools that are capable of controlling the characteristics of the generated text has become particularly appealing.

37 *Both authors contributed equally to this research.

Authors' addresses: Moreno La Quatra, moreno.laquatra@unikore.it, Kore University of Enna, Piazza dell'Università, Enna,
 Italy, 94100; Giuseppe Gallipoli, giuseppe.gallipoli@polito.it, Politecnico di Torino, Corso Duca degli Abruzzi 24, Turin,
 Italy, 10129; Luca Cagliero, luca.cagliero@polito.it, Politecnico di Torino, Corso Duca degli Abruzzi 24, Turin, Italy, 10129;

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee
 provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and
 the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored.
 Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires

- prior specific permission and/or a fee. Request permissions from permissions@acm.org.
- ⁴⁶ © 2024 Association for Computing Machinery.
- 47 2157-6904/2024/7-ART \$15.00
- 48 https://doi.org/XXXXXXXXXXXXXXX
- 49

30

31

1 2

3

Text Style Transfer (TST) is a well-known NLP task. It focuses on changing the style attributes of a piece of text from the source style to a target one (e.g., from an informal version to its formal one) while preserving the original message conveyed by the text. Changing the text style is relevant to a wide range of real-life applications ranging from online content moderation to intelligent writing assistants [13]. TST solutions may improve the user experience by enhancing the intelligibility and pertinence of the generated text as well as adapting the language to the current situation and writer/speaker's intent [9]. Importantly, style transfer must be achieved with minimal changes to the text to preserve the original content as much as possible.

In this work, we address TST in an unsupervised scenario, i.e., we assume that there is a lack of parallel annotated data to train sequence-to-sequence models [36]. The key challenges in unsupervised TST are (1) The preservation of the original content of the source text, and (2) The correct identification and replacement of the stylistic elements present in the textual content. In the absence of parallel training data, disentangling style and content is known to be particularly challenging [16]. On the other hand, unsupervised TST approaches are, broadly speaking, more resource-efficient as they do not involve labor-intensive training tasks [47].

We propose a new architecture for TST relying on Cycle-consistent Generative Adversarial Networks (CycleGANs). CycleGANs exploit the cycle-consistency principle for self-supervised adversarial learning. In the context of TST, they have recently emerged as promising sequence-to-sequence approaches to disentangle text style and content [12].

Existing CycleGAN-based approaches to TST face the following issues [12]:

- I1) **Self-supervision in the latent space**: they encode/decode the input/output text and employ fully-connected neural networks to implement the generator and discriminator models. This makes content and style information tightly connected to the text embedding representation, facing issues while coping with mixed-style content, i.e., input textual sequences that are partly in the original style (e.g., informal) and partly in the target one (e.g., formal).
- I2) Recurrent networks: generators and discriminators consist of LSTM networks, which are known to be suboptimal for coping with long-term text dependencies. Although the interest of the NLP community has already shifted towards the use of Transformer encoder-decoder architectures [38], to the best of our knowledge existing CycleGAN-based TST approaches do not rely on Transformers yet.
- I3) Weak enforcement of the target style in adversarial learning: since in the adversarial learning process the discriminator distinguishes between real and fake sentences without explicitly taking the style of the generated text into account, the target style is weakly enforced.

Our approach overcomes the limitations of existing approaches by introducing the following innovative features:

- Self-supervision at the sequence level: to overcome issue I1, it applies self-supervision, based on the cycle-consistency principle [47], directly to the raw input sequences. During the training process, the adversarial loss ensures that the generated text is indistinguishable from the target text, whereas the cycle-consistency loss ensures that the mapping between the source and target text styles is invertible.
- **CycleGANs using Transformers**: to overcome issue I2, it adopts a self-supervised approach based on CycleGANs [47] which automatically learns the mapping between the original and target styles without the need for paired data. The proposed framework consists of two generators and two discriminators. All of them are based on the Transformer architecture [38].

• **Classifier-guided text generation**: to overcome issue I3, the CycleGAN generators leverage a pre-trained classifier performing text style prediction. The classification loss returned by the classifier is integrated into the generators' loss functions to guide the text generation process. The style classifier is aimed to guide the generators to produce text with the desired style attributes while maintaining the original content's meaning.

The empirical results, achieved on benchmark TST datasets for sentiment and formality transfer, show the superior performance of the proposed approach:

- Against state-of-the-art unsupervised TST models: we compare the performance of our approach with that of recently proposed unsupervised approaches to TST, including Transformer-based and CycleGAN-based ones [12]. The presented architecture outperforms all the tested competitors, e.g., +6.8 points of SacreBLEU on the GYAFC dataset (see Tables 4 and 5).
 - On mixed-style inputs: we run extensive experiments on TST tests suited to a mixed-style scenario. The results, exemplified in Figure 1, confirm the superior performance of our approach while coping with mixed-style inputs.
 - Against Large Language Models: we also compare our approach with a state-of-the-art open-source Large Language Model with a similar number of parameters, i.e., Llama2-7B [37]. The results show that our approach averagely performs better on benchmark data and is more robust than the tested LLM to mixed-style inputs.
 - In a human evaluation: we carried out a human evaluation to qualitatively assess the quality of a sample of TST outcomes. The results are coherent with the quantitative performance metrics.

As an example, the results summarized in Figure 1 show that CycleGAN (our approach) generates output sequences that are most syntactically similar to the expected outcome (the higher ref-BLEU the better) on all the tested configurations of mixed-style inputs. The mixing ratio X%-Y% indicates the percentage ratio of original (X%) and target (Y%) style in the input. The performance of the Large Language Model (Llama2) is closer to that of CycleGAN when there is no mix (e.g., $X \approx 0\%$ or $Y \approx 0\%$), whereas is significantly worse in a mixed scenario (e.g., X=Y=50%).



Fig. 1. Performance comparison on mixed-style inputs. Dataset: GYAFC-music. Metric: ref-BLEU. Mixing
 ratios are in the form X%-Y%, where X and Y are the percentages of original and target style in the input,
 respectively. Approaches: CycleGAN (ours), Llama2-7B [37], CycleGAN latent (variant of CycleGAN without
 sequence-level cycle-consistency), CrossAlignment [34], and MultiDecoder [5].

In summary, the novelty of our TST approach lies in: (1) The application of cycle-consistency directly to the input sequence, making the approach more robust for content preservation, particularly when coping with mixed-style inputs (see the results in Section 5.9); (2) The adoption

of Transformers in a CycleGAN TST approach (see the empirical comparisons in Section 5.5); 148 and (3) The use of a style classifier to foster the generator to produce text in the target style (see 149 150 Section 4.1.1 for further details).

RELATED WORKS 152 2

153 According to a recently proposed categorization [9], existing TST methods can be classified as (1) 154 Parallel supervised, if they are trained on known pairs of text with different styles; (2) Non-parallel 155 supervised, if the style labels are available but the matching between text pairs is missing; (3) Purely 156 unsupervised, if the style labels are not available. 157

Parallel supervised or semi-supervised approaches (e.g., Shang et al. [33], Wang et al. [39], Xu et al. 158 [43]) require large-scale style-to-style parallel data, i.e., examples of parallel sentences conveying the 159 same message with different style attributes. However, their generation is extremely labor-intensive. 160 Conversely, non-parallel supervised approaches are trained on large text corpora annotated with style labels. Relaxing the constraint of having style-to-style text pairs makes the problem challenging 162 yet more tractable in real scenarios. This paper falls into the latter category. 163

Non-parallel supervised methods need to address the following issues:

- Content preservation: it involves maintaining the original textual content while transforming the text style. Preserving the underlying meaning, semantic information, and structural characteristics of the input text is essential to ensure the coherence and fidelity of the generated output. However, achieving effective content preservation while simultaneously changing the text style is a non-trivial task.
- *Style-content disentanglement*: it refers to the process of correctly separating the style attributes from the content in the text. This disentanglement is challenging because style and content are inherently intertwined and strongly related to each other. Modifying the style of a text without altering its content requires the model to accurately identify and manipulate the style-specific attributes while keeping the underlying content intact [16].

Style-content disentanglement can be achieved through different strategies:

- *Explicit disentanglement* [18, 42, 44]: it entails directly replacing the text with the original style attributes with new pieces of text that have the desired target style attribute. This approach explicitly separates the style and content. However, it can be applied only when style and content can easily be separated and the style transfer can be realized by changing only some selected words.
- Implicit disentanglement [5, 8, 26]: it learns two distinct latent representations, one for the content and the other for the style. By manipulating these separate representations, the model can ideally modify the style while preserving the content. Different techniques such as back-translation, attribute control generation and adversarial training are usually adopted to realize this approach.
- Without disentanglement [3, 7, 23]: the style-content separation is concealed and the model does not explicitly distinguish between them during the style transfer process. This approach aims at seamlessly transforming the style attributes while implicitly capturing and preserving the underlying content.

In our method, we adopt a strategy without disentanglement. Recent approaches to TST without disentanglement have explored the combination of linguistic graph structures and Transformerbased architectures [35]. An extensive review of existing techniques can be found in [9].

Adversarial learning has already been successfully employed to model style-content disentanglement and achieved fairly good content preservation. Recent works [1, 12, 21, 46] have already

195 196

151

161

164

165

166

167

168

169

170

171

172

173

174 175

176

177

178

179

180

181

182

183

184

185

186

187

188

189

190

191

192

193

adopted Generative Adversarial Networks (GANs) and cycle-consistency for non-parallel supervised Text Style Transfer. The key differences with the present work are summarized below.

- Zhao et al. [46] propose an encoder-decoder framework where text style and content are encoded into two distinct latent vectors (i.e., implicit disentanglement). The encoding and decoding functions are coupled with a style discrepancy loss, which models the style shift from the original domain to the target one, and with a cycle-consistency loss, which ensures content preservation. Unlike Zhao et al. [46], our approach adopts CycleGANs [47] and is without disentanglement.
 - Chen et al. [1] present a GAN framework that leverages optimal transport and uses the feature mover's distance [41] as training loss. Unlike the present work, they adopt the cycle-consistency principle only for addressing the task of unsupervised deciphering in the latent feature space, relying on LSTM networks for text generation and convolutional networks as sentence feature extractors.
 - Huang et al. [12] adopt CycleGANs by imposing the cycle-consistent constraint in the continuous latent space. They rely on the LSTM architecture to encode/decode the input/output sequence and employ a two-layer fully-connected neural network to implement the generator and discriminator models. In contrast, our proposed approach performs adversarial training on the raw text sequences and computes the cycle-consistency loss at the text level, allowing for more fine-grained control of the text attribute style.
 - Lorandi et al. [21] focus on sentiment transfer using CycleGANs and LSTMs. In contrast, our approach explores multiple style attributes, utilizes Transformer architectures, and integrates a style classifier to enhance style transfer quality and fidelity.

Recently, some research has explored the use of Large Language Models (LLMs) to address TST. For instance, Reif et al. [30] propose an augmented zero-shot learning strategy showing promising results on various TST tasks without requiring fine-tuning or exemplars in the target style. An empirical comparison between our method and an LLM can be found in Section 5.6.

3 PRELIMINARIES

Table 1. Summary of notations for the Text Style Transfer task.

Symbol	Description
<i>A</i> , <i>B</i>	Source/Target style
χ_s	Textual corpora with style s
x_s, y_s	Input (Output) sequence in style s
\mathcal{F}, \mathcal{G}	Mapping function from style $A(B)$ to $B(A)$
$G_{A \to B}, G_{B \to A}$	Generator from style $A(B)$ to $B(A)$
D_A, D_B	Discriminator for style $A(B)$
SC	Style classifier
L	Overall loss function
$\mathcal{L}_{G_{A \to B}}, \mathcal{L}_{G_{B \to A}}$	Generator loss
$\mathcal{L}_{G_{D_A}}, \mathcal{L}_{G_{D_P}}$	Adversarial loss
$\mathcal{L}_{cuc_{A} \rightarrow B \rightarrow A}, \mathcal{L}_{cuc_{B} \rightarrow A \rightarrow B}$	Cycle-consistency loss
$\mathcal{L}_{stule_A}, \mathcal{L}_{stule_B}$	Classifier-guided loss
$\mathcal{L}_{D_A}, \mathcal{L}_{D_B}$	Discriminator loss
$\mathcal{L}_{D_A}^{real}, \mathcal{L}_{D_A}^{fake}, \mathcal{L}_{D_B}^{real}, \mathcal{L}_{D_B}^{fake}$	Real/Fake sample discriminator loss
$\lambda_{gen}, \lambda_{cyc}, \lambda_{style}, \lambda_{dis}$	Loss scaling factors

In this section we introduce the preliminary concepts and formally state the problem under consideration. For the sake of readability, the notation used throughout the section is summarized in Table 1.

Text Style Transfer (TST). TST aims to learn a mapping function \mathcal{F} that transforms an input text x_A with source style A into its transferred version x_B with target style B. Similarly, function \mathcal{G} applies the reverse transformation, i.e., from x_B to x_A . Unlike style-conditioned text generation [14], in TST the transformation preserves the original content while transferring the style from A to B.

$$\mathcal{F}: x_A \to x_B \mid x_A \qquad \qquad \mathcal{G}: x_B \to x_A \mid x_B \tag{1}$$

Hereafter, we will consider the level of formality (i.e., formal or informal) or the sentiment score (i.e., positive or negative) as style attributes. The main TST complexity lies in the tight connection between content and style. For example, the level of formality of a piece of text is often determined not only by a particular linguistic register but also by other characteristics such as syntax and orthography. For the sake of simplicity, we also assume to be in a binary style transfer scenario¹.

Cycle-consistent Generative Adversarial Networks (CycleGANs). Our goal is to address TST by leveraging Cycle-consistent Generative Adversarial Networks (CycleGANs). They are a class of Generative Adversarial Networks (GANs) that can learn the mapping function between two domains without the need for parallel data [47]. Although they have been introduced in the field of Computer Vision, CycleGANs are general-purpose architectures that can be used to accomplish a variety of tasks, including TST. The use of CycleGANs enables the adoption of a self-supervised paradigm, relaxing the constraint on the availability of parallel textual data.

⁸ CycleGAN architectures typically comprise four models, including two generators and two
 ⁹ discriminators. The generators learn the mapping functions while the discriminators ensure the
 ⁰ quality of the generated outputs. In the following, we outline the general formulation of CycleGAN
 ¹ training objectives. For a detailed description of both the architecture and the training process
 ² specific to the task under consideration, please refer to Section 4.

Let X and Y be two domains with training examples $x_i \in X$ and $y_j \in Y$. The corresponding data distributions can be denoted as $x \sim p_{data}(x)$ and $y \sim p_{data}(y)$. The generators F and G aim to learn the following mappings: $\mathcal{F} : X \to Y$ and $\mathcal{G} : Y \to X$. The discriminator D_Y aims to distinguish between real samples y and generated samples F(x). Similarly, D_X discriminates between x and $G(y)^2$.

CycleGAN training involves two objectives: adversarial losses and cycle-consistency loss. Adversarial losses try to match the distribution of generated samples with the data distribution in the target domain. Specifically, for the mapping function $\mathcal{F} : X \to Y$, it can be expressed as follows:

$$\mathcal{L}_{\text{GAN}}\left(F, D_Y, X, Y\right) = \mathbb{E}_{y \sim p_{data}(y)}\left[\log D_Y(y)\right] + \mathbb{E}_{x \sim p_{data}(x)}\left[\log\left(1 - D_Y(F(x))\right)\right]$$
(2)

It is an adversarial loss since *F* aims to minimize it against an adversary D_Y that tries to maximize it. Similarly, it is possible to define an adversarial loss for the mapping function $\mathcal{G} : Y \to X$.

The cycle-consistency loss constrains the mapping functions to ensure that their sequential application to the input sample x allows for its reconstruction. Additionally, it addresses the mode collapse problem [47]. By combining the reconstruction constraint of both mapping functions, it can be defined as follows:

$$\mathcal{L}_{cyc}(F,G) = \mathbb{E}_{x \sim p_{data}(x)} \left[\|G(F(x)) - x\|_1 \right] + \mathbb{E}_{y \sim p_{data}(y)} \left[\|F(G(y)) - y\|_1 \right]$$
(3)

 $\frac{291}{^{1}\text{The multiple style transfer problem is out of the scope of the present work but, as discussed in [12], can be addressed by factorizing the problem into binary subtasks.}$

²With a slight abuse of notation, F(x) or $\mathcal{F}(x)$ will be used interchangeably hereafter for the sake of simplicity.

ACM Trans. Intell. Syst. Technol., Vol. 1, No. 1, Article . Publication date: July 2024.



Fig. 2. The figure illustrates the training process for cycle A to B to A; the training process for cycle B to A to 310 B is similar but the roles of the source and target texts are reversed. The generated text is reported using dashed lines and style A and style B are illustrated in blue and red, respectively. Numbered circles indicate the components of the loss function used to train the architecture. 313

Generator

This general formulation of CycleGAN training objectives can easily be adapted to the Text Style 315 Transfer task. The main difference lies in defining the two domains, X and Y, as the input and 316 target styles A and B. 317

METHOD 4 319

295

296

297

298

299 300

301

308 309

311

312

314

318

329

330

Figure 2 shows a sketch of the proposed method. The objective is to learn the mapping between 320 the two styles using two generators, $G_{A \to B}$ and $G_{B \to A}$, and two discriminators, D_A and D_B . These 321 components work together to learn the mapping between the source and target styles. A detailed 322 description of the generator and discriminator characteristics is given below. 323

In addition to the generators and discriminators, we also use an external, pre-trained style 324 classifier, hereafter denoted by SC. This model aims to classify the style of a given text sample. 325 During the training process of the CycleGAN model, the generators receive feedback from the 326 style classification model on the style of the generated content. This feedback is exploited by the 327 generators to effectively produce text pieces with the desired style attribute. 328

4.1 Generator

The purpose of the generators $G_{A\to B}$ and $G_{B\to A}$ is to learn the transformation between the 331 source and target pieces of texts. A modification of a specific text attribute, such as style or 332 sentiment, must preserve the original content. The generator $G_{A \rightarrow B}$ takes a sequence of tokens, 333 $x_A = (x_{A,1}, x_{A,2}, \dots, x_{A,N})$, as input, where $x_{A,n}$ is the *n*-th token in the sequence. The output of 334 the generator is a sequence of tokens, $y_B = (y_{B,1}, y_{B,2}, \dots, y_{B,M})$, where $y_{B,m}$ is the *m*-th token in 335 the output sequence. It is worth noting that the lengths of the input and output sequences may be 336 different (i.e., $N \neq M$). The generator $G_{B \rightarrow A}$ handles a similar process, transforming the piece of text 337 written in style B into a text written in style A. The key point is that the input and output sequences 338 are not paired, and the model cannot be trained based on prior knowledge of the expected output. 339

The proposed method involves two cycles, $A \to B \to A$ and $B \to A \to B$, which operate as 340 follows. For the cycle $A \rightarrow B \rightarrow A$, the generator $G_{A \rightarrow B}$ is trained to predict the output sequence 341 y_B given the input sequence x_A . The output of this generator is then fed to the discriminator D_B , 342

which aims at distinguishing between samples drawn from the original distribution and those generated by the generator $C_{\rm eff}$, which transform the input tout's style to the target style. The

generated by the generator $G_{A \to B}$, which transfers the input text's style to the target style. The output of the generator is also fed back to the generator $G_{B \to A}$, which transforms the style of the generated text back to the original style. The output of the second generator, $y_A = G_{B \to A}(y_B)$, corresponds to the reconstructed text that should be as close as possible to the original input text. The generators $G_{A \to B}$ and $G_{B \to A}$ are trained to minimize the following loss functions:

$$\mathcal{L}_{G_{A \to B}} = \lambda_{gen} \mathcal{L}_{G_{D_B}} + \lambda_{cyc} \mathcal{L}_{cyc_{A \to B \to A}} + \lambda_{style} \mathcal{L}_{style_B}$$
(4)

$$\mathcal{L}_{G_{B\to A}} = \lambda_{gen} \mathcal{L}_{G_{D_A}} + \lambda_{cyc} \mathcal{L}_{cyc_{B\to A\to B}} + \lambda_{style} \mathcal{L}_{style_A}$$
(5)

Here, $\mathcal{L}_{G_{D_R}}$ (illustrated in point 1 in Figure 2) and $\mathcal{L}_{G_{D_A}}$ are the adversarial losses (see Equation 2) and represent the feedback from the corresponding discriminator (i.e., the extent to which the generator is able to generate text that is indistinguishable from the target), whereas $\mathcal{L}_{cyc_{A \to B \to A}}$ (illustrated in point 3 in Figure 2) and $\mathcal{L}_{cyc_{B\to A\to B}}$ are the cycle-consistency losses (see Equation 3) computed at the end of the corresponding cycle by comparing the output of the second generator to the input sequence. More formal definitions of the discriminator and cycle losses are available in Sections 4.2 and 4.3, respectively. \mathcal{L}_{style_B} and \mathcal{L}_{style_A} are the style classifier losses that are computed using the pre-trained style classifier. These components of the loss function (represented in point 2 in Figure 2), aim at ensuring that the generator is able to generate text that is consistent with the target style. \mathcal{L}_{style_B} and \mathcal{L}_{style_A} corresponds to the binary cross-entropy loss between the predicted style and the target style (known according to the transformation being learned). The classifier-guided loss is computed using the pre-trained style classifier but only the generator is updated using this loss (e.g., the pre-trained style classifier is not trained during the adversarial training process). The classifier-guided loss can be formalized as follows:

$$\mathcal{L}_{style} = -\frac{1}{N} \sum_{i=1}^{N} [y_i \log SC(x_i) + (1 - y_i) \log(1 - SC(x_i))]$$
(6)

where *N* is the number of samples in the batch, x_i is the input sequence, y_i is the target label (i.e., 1 for style B and 0 for style A), and $SC(x_i)$ is the output of the style classifier (see Section 4.1.1 for further details) for the input sequence x_i classified using the [CLS] token. The loss is calculated by taking the average over all samples in the batch.

A different hyperparameter is associated with each of the three loss components in Equation 4 (and 5): specifically, λ_{gen} , λ_{cyc} and λ_{style} respectively control the relative importance of \mathcal{L}_{GD_B} (\mathcal{L}_{GD_A}) , $\mathcal{L}_{cyc_{A\to B\to A}}$ $(\mathcal{L}_{cyc_{B\to A\to B}})$ and \mathcal{L}_{style_B} (\mathcal{L}_{style_A}) .

For the cycles $A \rightarrow B \rightarrow A$ and $B \rightarrow A \rightarrow B$ the process is quite similar: the generators and discriminators operate to learn the transformation between the source and target styles.

4.1.1 Style classifier. The aim of the style classifier loss is to ensure the alignment with the target style, complementing content preservation and style transfer produced by the cycle consistency loss. It provides tailored guidance for accurate style transformation. Importantly, the discriminator, described in Section 4.2, distinguishes between real and fake sequences without taking the style of the generated output into account. Although it identifies out-of-distribution samples, the adversarial learning process weakly enforces the target style of the output text. To overcome this issue, the style classifier aims to provide explicit feedback to the generator on the style quality of the generated texts, thus mitigating the limitations of adversarial learning in TST.

4.2 Discriminator

The discriminators D_A and D_B are responsible for distinguishing between real and generated text. In line with the original GAN framework [6], the discriminators are trained to maximize the

9

probability of correctly classifying real and generated text. Specifically, D_A is trained to distinguish between the source texts and the output of the generator $G_{B\to A}$, where the source texts are samples drawn from the source distribution. Similarly, D_B is trained to distinguish between the target texts, which are samples drawn from the target distribution, and the output of the generator $G_{A\to B}$. The discriminators are trained to minimize the following loss functions:

$$\mathcal{L}_{D_A} = \mathcal{L}_{D_A}^{real} + \mathcal{L}_{D_A}^{fake} \tag{7}$$

$$\mathcal{L}_{D_B} = \mathcal{L}_{D_B}^{real} + \mathcal{L}_{D_B}^{fake} \tag{8}$$

where $\mathcal{L}_{D_A}^{real}$ and $\mathcal{L}_{D_B}^{real}$ denote the losses computed using data sampled from the source domain and the target domain, respectively, whereas $\mathcal{L}_{D_A}^{fake}$ and $\mathcal{L}_{D_B}^{fake}$ denote the losses computed using the output sequences of the generators $G_{B\to A}$ and $G_{A\to B}$, respectively. The weight of the discriminator losses in the overall objective function is controlled by a hyperparameter λ_{dis} .

Each term of the discriminator loss (i.e., $\mathcal{L}_{D_A}^{real}$, $\mathcal{L}_{D_A}^{fake}$, $\mathcal{L}_{D_B}^{real}$, and $\mathcal{L}_{D_B}^{fake}$) is defined as a Binary Cross-Entropy loss which, for a given discriminator D, is given by:

$$\mathcal{L}_{BCE} = -\frac{1}{N} \sum_{i=1}^{N} [y_i \log D(x_i) + (1 - y_i) \log(1 - D(x_i))]$$
(9)

where *N* is the number of samples in the batch, x_i is the input sequence, y_i is the target label (i.e., 1 for real text and 0 for generated text), and $D(x_i)$ is the output of the discriminator for the input sequence x_i classified using the [CLS] token. The loss is calculated by taking the average over all samples in the batch.

The adversarial losses computed using the output of the discriminators are back-propagated to the generators, allowing them to learn to generate text that is consistent with the data sampled from the target domain. By utilizing Transformer-based models as discriminators, we can efficiently learn the text style consistency within the target domain, thus improving the overall effectiveness of the training process.

4.3 Cycle consistency

The goal of the proposed method is to learn the mapping between the source and target domains. 423 In Figure 2 we illustrate the process for the case in which the source domain is A while B is the 424 target domain. However, the process is analogous the other way around. Given an input sequence 425 x_A in the source domain, the generator $G_{A \to B}$ is trained to generate a sequence y_B in the target 426 domain. However, due to the lack of parallel annotated data in the target domain during training, 427 the generator $G_{A \to B}$ is unable to directly learn the mapping between x_A and y_B . To address this 428 issue, the cyclic architecture first generates a sequence y_B in the target domain and then transforms 429 it back to the source domain using the generator $G_{B\to A}$. The output of such a generator, y_A , is then 430 compared to the input sequence x_A using a cycle-consistency loss. This loss is computed using the 431 cross-entropy loss between the output of the second generator and the input sequence and is used 432 to train the generator $G_{B \to A}$. 433

Specifically, each generator is a sequence-to-sequence model that is trained to minimize the cross-entropy loss between the generated sequence and the target sequence defined as follows:

$$\mathcal{L}_{cyc} = -\frac{1}{N \cdot T_{\text{total}}} \sum_{n=1}^{N} \sum_{t=1}^{T-1} y_{nt} \log(p_{t|t-1})$$
(10)

where *N* is the number of samples in the batch, T_{total} is the total number of tokens across all samples in the batch, *T* is the length of the sequence, y_{nt} is the target token at position *t* in the sequence *n*,

439

434

398 399

400 401

412

413

414

415

416

417

418

419

420 421

and $p_{t|t-1}$ is the probability of the token at position *t* given the previous tokens in the sequence. The loss is calculated by taking the average over all samples in the batch.

Given the self-supervised nature of our method, the target sequence is not available during the initial transformation from the source domain to the target domain (i.e., $A \rightarrow B$). The subsequent transformation from the target domain back to the source domain (i.e., $B \rightarrow A$) aims to reconstruct the original input sequence. At this stage, the expected output is the input sequence x_A , which can be used to compute the cycle-consistency loss. Therefore, the cycle-consistency loss is computed using both the output of the generator $G_{B\rightarrow A}(y_B) = y_A$ and the input sequence x_A (for the cycle $A \rightarrow B \rightarrow A$). A similar process occurs for the cycle $B \rightarrow A \rightarrow B$.

4.4 Objective function

The full objective function is a combination of various loss functions, with each component contributing to a specific aspect of the Text Style Transfer task. The loss functions include the generator loss, the cycle-consistency loss, the style loss, and the discriminator loss, each weighted by a hyperparameter λ . The final formulation can be expressed as follows:

$$\mathcal{L} (G_{A \to B}, G_{B \to A}, D_A, D_B) = \lambda_{gen} \mathcal{L}_{G_{D_B}} + \lambda_{cyc} \mathcal{L}_{cyc_{A \to B \to A}} + \lambda_{style} \mathcal{L}_{style_B} + \lambda_{gen} \mathcal{L}_{G_{D_A}} + \lambda_{cyc} \mathcal{L}_{cyc_{B \to A \to B}} + \lambda_{style} \mathcal{L}_{style_A} (11) + \lambda_{dis} \mathcal{L}_{D_A} + \lambda_{dis} \mathcal{L}_{D_B}$$

It includes the adversarial losses for both generators ($G_{A\rightarrow B}$ and $G_{B\rightarrow A}$), and discriminators (D_A and D_B), the cycle-consistency losses for both style transfer directions, and the style losses for both domains. Additionally, weighting factors (λ_{gen} , λ_{dis} , λ_{cyc} and λ_{style}) are used to balance the importance of each component in the overall objective.

It is worth noting that, while it is possible to have separate weighting factors for each direction, we employ identical weighting hyperparameters for both style transfer directions to maintain simplicity and minimize the complexity of configuration options. This choice allows us to avoid the need for justifying or making any prior assumption concerning distinct values for each direction. By ensuring uniformity in the weighting factors, we establish a balanced optimization process that treats both directions equally. The proposed implementation can easily be extended to accommodate different weighting factors for each style transfer direction if required by the specific use case. Finally, we formulate the overall optimization problem as follows:

$$G_{A \to B}^*, G_{B \to A}^* = \arg \min_{G_{A \to B}, G_{B \to A}} \max_{D_A, D_B} \mathcal{L} \left(G_{A \to B}, G_{B \to A}, D_A, D_B \right)$$
(12)

which expresses the min-max game played between each pair of generator-discriminator models [47].

4.5 Extension to Multiple Styles

The current approach is designed to handle only a specific pair of source and target styles. A 480 straightforward method to handle more than two styles is to train separate pairwise TST architec-481 tures. However, this leads to scalability issues. An alternative, more efficient solution is to prompt 482 the generator with specific instructions [2] indicating the desired style transformation. For example, 483 by using purposefully crafted prompt tokens like [A->B] for converting text from style A to style 484 B, or [A->C] for converting text from style A to style C, each generator can be trained to handle 485 multiple style conversions. This approach maintains the self-supervised nature of the architecture, 486 enabling generators to convert from any style to any other style. However, implementing this 487 method requires careful consideration of model training and architecture adjustments, which are 488 beyond the scope of the current work (see Section 6 for a discussion of the future research lines). 489

490

451

452

453

454

455

462

463

464

465

466

467

468

469

470

471

472

473 474 475

476

477 478

491 5 EXPERIMENTAL EVALUATION

We evaluate the performance of the proposed method and compare it against recent TST approaches
 on benchmark data. We also perform various ablation studies to evaluate the following aspects:
 cycle-consistency in the latent space, impact of the cycle-consistency loss coefficient, and effect of
 the pre-trained style classifier.

To foster the reproducibility of our results, the models and code used for the implementation of the proposed framework are publicly available, for research purposes only, at https://github.com/gallipoligiuseppe/TST-CycleGAN under the license CC BY-NC-SA.

5.1 Datasets

We consider three benchmark datasets related to two different TST tasks, i.e., sentiment transfer and formality transfer.

Sentiment transfer. The Yelp dataset [34] collects restaurant reviews. Based on their rating,
 reviews are labeled as positive or negative (a rating of 4 or 5 corresponds to a positive label, whereas
 a rating below 3 is negative). The dataset includes a test set with four human references per sentence.
 Train and validation sets are suited to non-parallel supervised TST as they are annotated with style
 attributes but the matching between text pairs is missing. For the sake of reproducibility, we use
 the same train/validation/test splits as in Li et al. [18] (see Table 2).

	# train	# validation	# test
negative	177,218	2,000	500
positive	266,041	2,000	500
total	443,259	4,000	1,000

Table 2. Yelp dataset statistics.

Formality transfer. Grammarly's Yahoo Answers Formality Corpus (GYAFC) [28] is a parallel dataset consisting of informal-to-formal sentence pairs. It comprises sentences from two different domains, i.e., Family & Relationships (family, in short) and Music & Entertainment (music, in short). Although the dataset includes parallel sentences, to simulate the scenario of self-supervised style transfer we select only the source sentences from the train set. The validation and test sets, on the other hand, include annotated sentences for both domains and are used to evaluate the performance of the proposed method (see Table 3).

Table 3. GYAFC dataset statistics.

		# train	# validation	# test
family	informal	51,967	2,788	1,332
	formal	51,967	2,247	1,019
	total	103,934	5,035	2,351
	informal	52,595	2,877	1,416
music	formal	52,595	2,356	1,082
	total	105,190	5,233	2,498

540 5.2 Metrics

We evaluate our model using a suite of established evaluation metrics [7, 23, 31]. Specifically, to
 quantify content preservation we compute the SacreBLEU score [25] between the system outputs
 and the four human references³.

544 To evaluate the effectiveness of our approach for Text Style Transfer, we fine-tune a BERT-base 545 binary classifier [4] to compute the style accuracy metric. To distinguish it from the style classifier 546 used during model training, hereafter we will denote it by the oracle classifier. Its purpose is to 547 accurately classify the style of the input text and provide a reliable evaluation metric for the quality 548 of the generated text's style transfer. On the analyzed datasets, the oracle classifier respectively 549 achieves the accuracies of 98.5% (Yelp), 94.0% (GYAFC-family), and 94.6% (GYAFC-music). To 550 compute the style accuracy, according to prior works we also train a TextCNN [15] as oracle 551 classifier (beyond the BERT-base classifier). Its classification performance is, in general, satisfactory 552 on all the tested datasets (96.5% on Yelp, 93.2% on GYAFC-family, 93.8% on GYAFC-music) and 553 comparable to that of the BERT-base model. Both BERT-base and TextCNN classifiers were trained 554 on the same TST datasets under analysis: they were trained and tested on the corresponding 555 training and test splits, respectively. 556

Finally, we also provide a comprehensive performance score by computing the geometric mean (GM) and harmonic mean (HM) of the SacreBLEU and style accuracy scores.

5.3 Configuration settings

560 We implemented the proposed architecture using the Hugging Face Transformers library [40]. 561 As reference models for the generators and discriminators, we consider BART-base [17] (140M 562 parameters) and DistilBERT [32] (66M parameters), respectively. For the Yelp dataset we use the 563 case-insensitive variant of DistilBERT, whereas for GYAFC we use the case-sensitive version as 564 the input data contains case-sensitive text. We also run experiments with larger generator models 565 prioritizing the use of more powerful models for the most challenging (and resource-demanding) text 566 generation task. Specifically, for the generators we also consider BART-large (400M parameters) and 567 T5 [27] (with 60M, 220M, and 770M parameters for the small, base, and large versions, respectively). 568 As style classifier, we use a fine-tuned BERT-base (110M parameters) model. Note that although we 569 use the same model as for the oracle classifier, this is not necessarily the case. The proposed TST 570 architecture can be trained end-to-end, enabling the simultaneous learning of the style transfer 571 functions in both directions. 572

We use the validation SacreBLEU score as the reference metric to identify the best-performing training configurations and the optimal model checkpoints. Then, we evaluate them on the test set. The SacreBLEU score is calculated separately for each TST direction and the average of these two values is used as an overall score. Note that for the Yelp dataset human references are not available for the validation set. To overcome this issue, we optimize the geometric mean of the SacreBLEU score calculated between the system outputs and the source sentences, and the style accuracy.

Training details. Similar to [12], we train the model in a self-supervised setting for both tasks, even though the GYAFC dataset is a parallel corpus. Thus, for our purposes, the alignments between sentence pairs are neglected.

Based on our preliminary experiments, we observed that the impact of the hyperparameters λ_{gen} , λ_{dis} and λ_{style} is negligible. Therefore, for the sake of simplicity, hereafter we will set $\lambda_{gen} = \lambda_{dis} =$ $\lambda_{style} = 1$. We tune the following two hyperparameters: the learning rate and the loss scaling factor λ_{cyc} . The learning rate, which controls the magnitude of the weight updates during training, is

12

557

558 559

573

574

575

576

577

578

579

580

581

586

³We adopt the implementation of the sacrebleu metric available at https://github.com/mjpost/sacreBLEU.

ACM Trans. Intell. Syst. Technol., Vol. 1, No. 1, Article . Publication date: July 2024.

⁵⁸⁹ updated using a linear scheduler, which linearly decreases the learning rate from a maximum value, ⁵⁹⁰ as reported in Appendix, to zero during the training process. Meanwhile, the λ_{cyc} hyperparameter ⁵⁹¹ controls the weight of the cycle-consistency loss in the overall objective function.

Given the computational demands of training such models and to reduce the number of con-592 figurations to be tested, we explore the hyperparameter space by considering values in the range 593 $[10^{-5}, 10^{-3}]$ for the learning rate and $\{0.1, 1, 10\}$ for the loss scaling factor λ_{cyc} . The optimal hyper-594 parameter values used throughout the experiments are reported in Appendix. It is worth noticing 595 that the selection of appropriate hyperparameters may affect the performance of the model. More-596 over, these hyperparameters were found to be optimal for the specific datasets and models used in 597 our experiments, and they may not necessarily generalize to other datasets or models. The optimal 598 values were determined through a combination of manual tuning and grid search, by evaluating 599 the model's performance of various hyperparameter combinations on the validation sets. 600

To balance computational efficiency and model performance, we set the maximum input sequence length to 64 since the average number of tokens is 8.88 ± 3.64 and 10.68 ± 4.12 for the Yelp and GYAFC datasets, respectively, and the batch size to 128 for BART-base, 32 for BART-large, 128 for T5-small, 64 for T5-base, and 8 for T5-large. We employ the AdamW optimizer [22] with β_1 0.9, β_2 0.999, and weight decay 10^{-2} .

To ensure consistent experimental conditions and hardware utilization, we utilize a single NVIDIA[®] V100 GPU with 32 GB of VRAM for both training and inference of all models.

5.4 Baselines

606

607 608

609

617

618

619

620

621

622

623

624

625

626

627

628

629

630

631

632

Existing unsupervised TST methods. We test the following methods: RetrieveOnly, DeleteOnly,
 DeleteAndRetrieve, and TemplateBased [18], BackTranslation [26], StyleEmbedding and Multi Decoder [5], CrossAlignment [34], UnpairedTranslation [42], UnsupervisedMT [45], DualRL [23],
 NASTLatentLearn [11], DeepLatent [7], ReinfRewards [31], MixAndMatch [24], MultiClass [3],
 FineGrainedST [19], LevenshteinEdit [29], GTAE [35], CycleAutoEncoder [12], and TextGANPG
 [21]⁴. Notice that we disregard existing supervised approaches to formality transfer because we
 deem their comparison with unsupervised methods unfair.

Cycle-consistency in the latent space. A key property of our approach is that it performs style transfer directly at the sequence level. Conversely, previous CycleGAN-based TST approaches apply transformations on the latent space. To evaluate our method's effectiveness against this approach, we explore two variants that conduct style transfer in the latent space. In these variants, we leverage the embedding space for style transfer. We decompose the generator network into encoder *E* and decoder *D* components, introducing two baseline models:

- (1) **Sentence-level**: this approach focuses on aligning representations generated by the encoders E_A and E_B with each other. Considering the case $A \rightarrow B \rightarrow A$, this is achieved by minimizing the L1 loss between the embeddings of the input sequence encoded by E_A and its corresponding version, predicted in the target style, and then encoded by E_B . To obtain sentence representations from token representations we use average pooling. The rationale behind this approach is to ensure that the semantic meaning of the input sequence is preserved while transferring its style.
 - (2) **Token-level**: in this baseline model, the focus is on preserving the content of the text at the token level by maximizing the similarity between the original input and the reconstructed

 ⁴To ensure a fair comparison, we recompute the results of the baseline methods using our own evaluation scripts. When the baseline methods produce lowercase outputs, we lowercase the human references and retrain the style classifiers on the lowercase versions of the datasets. Lower-cased *oracle classifiers* accuracies: 90.0% (GYAFC-family) and 91.1% (GYAFC-music);
 TextCNN models: 89.2% (GYAFC-family) and 88.8% (GYAFC-music).

641		ref-B avg	acc _{BERT}	GM	HM	acc _{CNN}	GM	HM
642	RetrieveOnly [18]	4.4	92.9	20.2	8.4	47.7	14.5	8.1
643	BackTranslation [26]	6.1	44.4	16.5	10.7	48.1	17.1	10.8
644	StyleEmbedding [5]	12.1	30.6	19.2	17.3	43.1	22.8	18.9
645	MultiDecoder [5]	16.1	25.6	20.3	19.8	42.8	26.3	23.4
646	CrossAlignment [34]	8.1	61.6	22.3	14.3	46.2	19.3	13.8
647	UnpairedTranslation [42]	5.2	58.1	17.4	9.5	47.2	15.7	9.4
648	DeleteOnly [18]	27.9	28.0	27.9	27.9	47.2	36.3	35.1
649	DeleteAndRetrieve [18]	21.0	52.4	33.2	30.0	45.4	30.9	28.7
650	TemplateBased [18]	31.9	39.4	35.4	35.2	46.0	38.3	37.7
651	UnsupervisedMT [45]	30.6	65.1	44.6	41.6	45.3	37.2	36.5
652	DualRL [23]	36.6	53.2	44.1	43.3	41.7	39.1	38.9
653	NASTLatentLearn [11]	38.6	49.3	43.6	43.3	43.1	40.8	40.7
654	CycleGAN BART (base)	43.7	50.7	47.1	46.9	49.4	46.5	46.4
655	CycleGAN BART (large)	43.5	50.8	47.0	46.9	49.9	46.6	46.5
656	CycleGAN T5 (small)	42.1	38.7	40.4	40.3	39.6	40.8	40.8
657	CycleGAN T5 (base)	44.0	47.7	45.8	45.8	46.2	45.1	45.1
658	CycleGAN T5 (large)	45.4	59.5	52.0	51.5	58.1	51.4	51.0

Table 4. Results on the GYAFC-family dataset – ref-BLEU (ref-B avg), style accuracy with BERT and TextCNN (acc_{BERT}, acc_{CNN}), geometric mean (GM) and harmonic mean (HM) of ref-BLEU and style accuracy.

output. To achieve this, it minimizes the L1 loss between the token embeddings of the input sequence and its reconstructions. The token-level embeddings for the input sequence are obtained immediately after tokenization, before being fed into the model. The reconstructed tokens are taken from the output of the decoder after completing the cycle (i.e., $A \rightarrow B \rightarrow A$) in the CycleGAN architecture.

To assess the performance of the two described latent-based variants of our approach, we present experimental results in Section 5.7.

5.5 Evaluation and Comparison

Here we summarize the main results of the empirical evaluations and performance comparisons separately for each style transfer domain. We also conduct a qualitative analysis of the generated outputs whose results are provided in Appendix.

Formality transfer. Tables 4 and 5 report the performance of our method variants (denoted by 673 the prefix name CycleGAN) and the baselines on the family and music domains of the GYAFC 674 dataset, respectively. The music domain has been shown to be more challenging and, in general, 675 less explored by previous TST studies than the family one. In both domains, our approach based 676 on T5 large performs best in terms of SacreBLEU scores compared to all the tested prior works. 677 More specifically, CycleGAN outperforms the other approaches in terms of ref-BLEU score (+6.8 vs. 678 the best-performing competitor), showing a higher capability of content preservation and a better 679 fluency of the generated text. Conversely, models achieving the highest accuracy scores significantly 680 perturb the original content as the corresponding ref-BLEU scores are fairly low, resulting in a less 681 faithful reproduction of the original meaning. Instead, the proposed approach achieves the best 682 balance between content preservation and style transfer. Among the tested CycleGAN variants, 683 those relying on larger generator models produce, as expected, consistently better ref-BLEU results 684 than the other ones. 685

ACM Trans. Intell. Syst. Technol., Vol. 1, No. 1, Article . Publication date: July 2024.

640

659 660

661

662

663

664

665

666

667 668

669

670

671

702

703

704

705

706

707

708

709

710

711

712

717

735

687	Table 5. Results on the GYAFC-music dataset - ref-BLEU (ref-B avg), style accuracy with BERT and TextCNN
688	(acc _{BERT} , acc _{CNN}), geometric mean (GM) and harmonic mean (HM) of ref-BLEU and style accuracy. * denotes
680	results from the paper.

	ref-B avg	acc _{BERT}	GM	HM	acc _{CNN}	GM	HM
DeepLatent [7]	21.4	77.8	40.8	33.6	56.2	34.7	31.0
ReinfRewards* [31]	29.2	-	-	-	-	-	-
MixAndMatch* [24]	27.7	-	-	-	-	-	-
CycleGAN BART (base)	43.6	57.2	49.9	49.5	57.8	50.2	49.7
CycleGAN BART (large)	42.0	43.1	42.5	42.5	43.8	42.9	42.9
CycleGAN T5 (small)	40.6	37.9	39.2	39.2	39.0	39.8	39.8
CycleGAN T5 (base)	42.0	45.4	43.7	43.6	47.7	44.8	44.7
CycleGAN T5 (large)	45.6	70.5	56.7	55.4	70.1	56.5	55.3

More detailed results on the most common formality transfer case, i.e., from informal to formal style, are given in Appendix. The results confirm the superior performance of CycleGAN T5 large compared to all the other methods (e.g., ref-BLEU +34.1 against ReinfRewards on GYAFC-music).

Sentiment transfer. Table 6 reports the results of our method variants (denoted by the prefix name *CycleGAN*) and the baselines on the Yelp dataset. Our method shows performance superior to all the other methods in terms of average SacreBLEU metric using the BART large model (CycleGAN 56.5 vs. 54.9 of the best-performing competitor). The better ability to preserve the original content is partly mitigated by the lower style accuracy which is, however, less critical for the sentiment transfer task (e.g., CycleGAN \approx 75% accuracy in sentiment transfer vs. \approx 50% in formality transfer). In fact, sentiment transfer commonly requires minimal modifications of the text to change its polarity.

In general, we claim that our model is able to achieve better content preservation thanks to the
 cycle-consistent structure of our architecture which is instrumental in preventing inappropriate or
 unnecessary modifications of the input text.

5.6 Formality and Sentiment transfer with Large Language Models

718 We perform an empirical comparison between our approach and a state-of-the-art open-source 719 Large Language Model, i.e., Llama2 model [37]. Specifically, we consider the 7B version to ensure 720 a fair comparison in terms of model size with the proposed architecture⁵. We employ it in both 721 zero-shot and few-shot settings: in the latter case, we experiment with varying number of examples 722 $k \in \{1, 3, 5, 10\}$ provided as input to the model. Few-shot examples consist of both the input sentence 723 and the corresponding expected output in the target style. Examples are randomly selected from 724 the parallel training sets for formality transfer datasets, whereas in the case of sentiment transfer, 725 since no parallel data is available, we manually annotate the expected outputs for the selected 726 examples. 727

- Based on preliminary experiments, we set the model's temperature hyperparameter at 0.6. We provide the LLM with the following prompt:
- 730 Transform the following sentence from [SRC] style to [TGT] style.
- 731 Apply only minimal changes and preserve the meaning of the sentence.
- Here you can find some examples of sentences in [SRC] style and corresponding sentences in [TGT] style:

⁷³⁴ ⁵Due to computational constraints, we employ a 16-bit quantization.

ACM Trans. Intell. Syst. Technol., Vol. 1, No. 1, Article . Publication date: July 2024.

739

736	Table 6. Results on the Yelp dataset - ref-BLEU (ref-B avg), style accuracy with BERT and TextCNN (acc _{BERT} ,
737	acc _{CNN}), geometric mean (GM) and harmonic mean (HM) of ref-BLEU and style accuracy. * denotes results
738	from the paper.

740		ref-B avg	acc _{BERT}	GM	HM	acc _{CNN}	GM	HM
741	RetrieveOnly [18]	6.4	99.9	25.3	12.0	93.0	24.4	12.0
742	BackTranslation [26]	10.5	94.6	31.5	18.9	95.3	31.6	18.9
743	StyleEmbedding [5]	42.2	7.9	18.3	13.3	8.8	19.3	14.6
744	MultiDecoder [5]	29.1	46.8	36.9	35.9	50.1	38.2	36.8
745	CrossAlignment [34]	20.0	72.8	38.2	31.4	74.9	38.7	31.6
746	UnpairedTranslation [42]	33.1	50.4	40.8	39.9	51.4	41.2	40.3
747	DeleteOnly [18]	29.1	85.3	49.8	43.4	87.9	50.6	43.7
748	DeleteAndRetrieve [18]	30.0	89.9	51.9	44.9	90.6	52.1	45.1
749	TemplateBased [18]	39.7	83.6	57.6	53.8	85.3	58.2	54.2
750	UnsupervisedMT [45]	41.3	95.7	62.9	57.7	97.4	63.4	58.0
751	MultiClass [3]	51.8	86.1	66.8	64.7	87.2	67.2	65.0
752	DualRL [23]	51.5	88.5	67.5	65.1	90.1	68.1	65.5
753	DeepLatent [7]	40.4	83.8	58.2	54.5	86.2	59.0	55.0
754	FineGrainedST [19]	16.2	91.4	38.5	27.5	91.5	38.5	27.5
755	LevenshteinEdit [29]	48.9	87.2	65.3	62.7	82.9	63.6	61.5
756	GTAE [35]	51.1	86.7	66.5	64.3	85.9	66.2	64.1
757	NASTLatentLearn [11]	54.9	78.4	65.6	64.6	81.8	67.0	65.7
758	MixAndMatch [24]	46.6	88.3	64.1	61.0	81.7	61.7	59.3
759	CycleAutoEncoder* [12]	22.5	-	-	-	86.9	44.2	35.7
760	TextGANPG* [21]	32.4	68.0	46.9	43.9	-	-	-
761	CycleGAN BART (base)	55.7	78.8	66.3	65.3	77.8	65.8	64.9
762	CycleGAN BART (large)	56.5	75.1	65.1	64.5	74.6	64.9	64.3
763	CycleGAN T5 (small)	53.0	78.0	64.3	63.1	78.2	64.4	63.2
764	CycleGAN T5 (base)	54.2	76.6	64.4	63.5	77.3	64.7	63.7
765	CycleGAN T5 (large)	55.3	72.9	63.5	62.9	73.7	63.8	63.2

```
767 Input ([SRC] style): [SRC_EXi]
```

```
768 Output ([TGT] style): [TGT_EXi]
```

```
769 .
```

766

...
Input ([SRC] style): [SRC_INPUT]

770 Input ([SRC] style): [SRC 771 Output ([TGT] style):

where we replace [SRC] and [TGT] with the actual source and target styles, [SRC_EXi] and [TGT_EXi] with the source and target sentences for each of the *k* examples (for k > 0), and [SRC_INPUT] with the current test sample to be transferred.

Table 7 reports the results achieved for both formality and sentiment transfer tasks, while more 775 776 detailed results for the informal-to-formal transfer can be found in Appendix. In both tasks, the ref-BLEU and accuracy performance generally increases while providing more input examples until 777 reaching a steady state. This is probably due to the fact that, when providing numerous examples, 778 some noise may be introduced, potentially misleading the model. Surprisingly, the ref-BLEU results 779 on the Yelp dataset for k = 1, 3 are worse than in the zero-shot setting. One possible explanation is 780 that, since in the sentiment transfer task style can often be transferred by modifying only a few 781 words, in the zero-shot setting the model may tend to apply fewer modifications, resulting in a 782 higher ref-BLEU score but lower accuracy. In the 1- or 3-shot settings, the accuracy increases at 783

784

785Table 7. Results on the GYAFC-family, GYAFC-music, and Yelp datasets of Llama2-7B-Chat model for varying786number of examples k in the 0/few-shot setting – ref-BLEU (ref-B avg), style accuracy with BERT and TextCNN787(acc_BERT, acc_CNN), geometric mean (GM) and harmonic mean (HM) of ref-BLEU and style accuracy. \star denotes788the results of our best models.

700									
789	dataset	k	ref-B avg	acc _{BERT}	GM	HM	acc _{CNN}	GM	HM
790		0	24.0	80.4	43.9	36.9	77.0	43.0	36.6
791		1	21.9	85.2	43.2	34.9	82.8	42.5	34.6
792		3	36.8	87.5	56.7	51.8	86.9	56.5	517
793	GYAFC-family	5	37.2	84.6	56.1	517	83.5	55.7	51.4
794		10	34.0	883	54.8	<i>J</i> 1.7	87.5	54.5	10.0
795		10	34.0	00.3	54.0	49.1	07.J	54.5	49.0
796		*	45.4	59.5	52.0	51.5	58.1	51.4	51.0
797		0	26.0	78.1	45.1	39.0	73.3	43.6	38.4
798		1	34.0	95.0	56.8	50.1	91.8	55.8	49.6
700	OVATO	3	40.1	94.1	61.5	56.3	91.9	60.7	55.8
800	GIAFC-music	5	37.4	92.2	58.7	53.2	90.4	58.1	52.9
800		10	38.3	92.1	59.4	54.1	90.2	58.8	53.7
801		*	45.6	70.5	56.7	55.4	70.1	56.5	55.3
802		0	42.7	83.2	59.6	56.5	79.6	58.3	55.6
803		1	34.0	92.4	56 1	497	88.9	55.0	492
804		2	26.1	02.0	57.0	52.0	80.4	E 6 0	E1 /
805	Yelp	3	30.1	92.9	57.9	52.0	89.4	50.8	51.4
806	1	5	43.1	91.2	62.7	58.5	87.4	61.4	57.7
807		10	53.4	84.8	67.3	65.5	82.9	66.5	64.9
808		*	56.5	75.1	65.1	64.5	74.6	64.9	64.3
111/11			1				1		

the expense of a lower ref-BLEU score. This is likely because the model requires more examples toadhere to the requirement of applying only minimal changes to the input sentences.

By comparing LLM results with those of our best models, we can state that our method consis-812 tently outperforms Llama2 in terms of content preservation on both tasks (i.e., +8.2 on GYAFC-family, 813 +5.5 on GYAFC-music, +3.1 on Yelp). In contrast, Llama2 achieves the highest accuracy scores, even 814 when compared to the other baselines in the formality transfer task. The results highlight that 815 the TST performance of Llama2 is fair without ad hoc fine-tuning. It is also worth noting that 816 model fine-tuning would require parallel data and thus is out of the scope of the present work. 817 In conclusion, our proposed approach confirms its superior performance in content preservation, 818 even when compared to a larger and more powerful Large Language Model. 819

5.7 Cycle-consistency in the Latent Space: Sentence-Level vs. Token-Level

In this section, we present the results of an ablation study conducted to compare the performance of the latent-based versions of our approach (described in Section 5.4). The purpose is to quantitatively compare these model variants with the proposed solution, which enforces the cycle-consistency constraint directly to the raw input sequence. To better isolate the effect of the cycle-consistency loss, we exclude the pre-trained style classifier and its corresponding loss term. Additionally, to ensure a fair comparison, we use the same generator model that achieved the best results in the corresponding task (i.e., T5 large for formality transfer and BART large for sentiment transfer).

The results on both tasks are shown in Table 8. As can be seen, the non-latent version (i.e., applying cycle-consistency on the raw input sequence) significantly outperforms both latent-based versions in terms of geometric mean and harmonic mean on both tasks. Upon closer inspection, in the formality transfer task, our approach achieves the best ref-BLEU scores, exhibiting an

833

809

820

Table 8. Ablation study. Results on the GYAFC-family, GYAFC-music, and Yelp datasets of latent-based 834 cycle-consistency losses - ref-BLEU (ref-B avg), style accuracy with BERT and TextCNN (acc_{BERT}, acc_{CNN}), 835 geometric mean (GM) and harmonic mean (HM) of ref-BLEU and style accuracy. 836

837										
020	dataset	model	latent	ref-B avg	acc_{BERT}	GM	HM	acc _{CNN}	GM	HM
030			-	44.4	49.2	46.7	46.7	47.9	46.1	46.1
839	GYAFC-family	CycleGAN T5 (large)	sentence	8.0	28.9	15.2	12.5	33.7	16.4	12.9
840			token	5.9	43.3	16.0	10.4	43.5	16.0	10.4
841			-	43.3	43.0	43.1	43.1	45.1	44.2	44.2
842	GYAFC-music	CycleGAN T5 (large)	sentence	8.1	23.1	13.7	12.0	36.3	17.2	13.2
843			token	6.3	59.1	19.3	11.4	59.1	19.3	11.4
045			-	56.9	73.9	64.8	64.3	73.1	64.5	64.0
844	Yelp	CycleGAN BART (large)	sentence	58.6	1.6	9.7	3.1	3.6	14.5	6.8
845	-		token	0.7	99.7	8.4	1.4	98.6	8.3	1.4

improvement of more than +35.0 points on both domains. Considering style accuracy, the latent 848 token-level version performs the best on the GYAFC-music dataset. However, it must be noted that 849 the corresponding ref-BLEU score is extremely low. In the sentiment transfer task, the sentence-850 level and token-level versions achieve the best results in ref-BLEU and style accuracy, respectively. 851 Nonetheless, they show remarkably low results in the other metric of interest (i.e., sentence-level 852 accuracy=1.6, token-level ref-BLEU=0.7). 853

After manually inspecting the generated outputs, we observed that in the sentence-level version, 854 in most cases, the input is simply copied to the output. The loss function used to train the model 855 aims at minimizing the discrepancy between the embedding of the input and transferred sentence, 856 therefore preserving the meaning. However, in the formality transfer task, where the output often 857 contains multiple copies of the input, there is a low overlap with the target sentence. In contrast, in 858 the sentiment transfer task, the input is copied to output without repetitions. Given that sentiment 859 transfer typically involves modifying only a few words, the sentence-level version achieves a high 860 ref-BLEU score. Considering accuracy scores, since the input text is not modified in the sentence-861 level version, its performance is low, especially in the sentiment transfer task (i.e., accuracy=1.6). 862 For the token-level version, the outputs are degenerate, i.e. the model almost generates the same 863 sentence. Consequently, the ref-BLEU scores are particularly low (e.g., 0.7 on the Yelp dataset), 864 while the accuracy is often very high (e.g., 99.7 on the Yelp dataset) if the (degenerate) sentences 865 are classified as belonging to the target style. 866

Overall, these results demonstrate the significant advantage achieved by directly enforcing 867 cycle-consistency constraints to the raw sequence, highlighting one of the main contributions of 868 our work. 869

Human Evaluation 5.8

Similar to [12, 23], we conducted a human evaluation to get a qualitative feedback on the TST 872 results. We recruited 12 volunteers, each of them meets the following criteria: she/he holds an 873 MSc or PhD degree, is proficient in English, and has a sufficient background in the Text Style 874 Transfer task. We randomly picked 50 test samples per dataset and style transfer direction (300 875 samples overall). For each source sample and target style, annotators were asked to evaluate the 876 quality of outputs generated by different systems. The outputs were presented in random order and 877 without disclosing the model each output was generated from. Specifically, for each task and dataset, 878 annotators evaluated the outputs produced by the following models: CycleGAN (ours), CycleGAN 879 latent (i.e., the latent sentence-based version of CycleGAN), Llama2, and the two corresponding 880 best baselines. 881

882

870

871

883	Table 9. Human Evaluation results on the GYAFC-family, GYAFC-music, and Yelp datasets - style accuracy
884	(Style), content preservation (Content), fluency (Fluency), average ratings (Avg) and success rate (Success).
885	* denotes scores for which $p < 0.05$.

886	dataset	model	Style	Content	Fluencv	Avg	Success
887		DualRL [23]	2.2*	3.0*	2.6*	2.6*	8.5%*
880		NASTLatentLearn [11]	2.3^{*}	2.9^{*}	2.7^{*}	2.6^{*}	$5.5\%^{*}$
890	GYAFC-family	CycleGAN latent T5 (large)	1.5^{*}	2.1^{*}	1.2^{*}	1.6^{*}	$0.5\%^*$
801		Llama2-7B-Chat	4.2^{*}	4.6^{*}	4.7	4.5^{*}	73.0%*
892		CycleGAN T5 (large)	3.5	4.9	4.7	4.4	55.0%
893		DeepLatent [7]	2.6^{*}	1.9*	2.5^{*}	2.3^{*}	$35.0\%^{*}$
894	CVAEC music	CycleGAN latent T5 (large)	1.7^{*}	2.6^{*}	1.4^{*}	1.9^{*}	$0.5\%^{*}$
895	GIARC-IIIusic	Llama2-7B-Chat	4.3 *	4.0^{*}	4.3	4.2^{*}	$63.0\%^*$
896		CycleGAN T5 (large)	4.1	4.8	4.5	4.5	74.4%
897		DualRL [23]	3.1*	3.2^{*}	3.5^{*}	3.3*	$26.0\%^{*}$
898		NASTLatentLearn [11]	2.9^{*}	3.2^{*}	3.1^{*}	3.1^{*}	$15.0\%^*$
899	Yelp	CycleGAN latent BART (large)	1.1^{*}	3.3*	4.2^{*}	2.9*	$0.5\%^*$
900		Llama2-7B-Chat	4.1	4.0^{*}	4.4^{*}	4.2^{*}	$64.5\%^{*}$
901		CycleGAN BART (large)	4.3	4.4	4.5	4.4	80.0%

The output sentences were evaluated using a 5-point Likert scale based on three criteria: (1) Style 904 accuracy, measuring the extent to which the generated sentence aligns with the target style; (2) 905 Content preservation, assessing how effectively the content of the input sentence is preserved; and 906 (3) Fluency, considering the overall fluency and linguistic correctness of the output text. Similar 907 to [23] and [18], we also calculate the average across the three criteria and denote a generated 908 output as "successful" if it receives a rating of 4 or 5 on all three criteria. 900

Table 9 reports the results achieved for both tasks, including a t-test for statistical significance. 910 In the formality transfer task, our approach excels in content preservation and fluency, achieving 911 the best performance. Moreover, it yields the highest average score and success rate on the GYAFC-912 music dataset. Conversely, Llama2 demonstrates the highest style transfer score for both domains, 913 and excels in terms of average score and success rate on the GYAFC-family dataset. Notably, our 914 approach and Llama2 outperform other systems across all metrics by a substantial margin (e.g., +2.9 915 and +2.0 on content preservation and style accuracy, respectively), especially the latent sentence-916 based version of our approach which exhibits the lowest performance. In the sentiment transfer task, 917 our model outperforms all baselines, including Llama2, on all metrics, thus confirming the superior 918 quality of the generated outputs. Broadly speaking, the human evaluations are mostly aligned with 919 the quantitative results and confirm the superior performance of our approach, particularly on 920 content preservation. Notably, we achieved exceptionally high scores in the formality transfer task 921 (i.e., 4.9 and 4.8 on the family and music domains, respectively), highlighting its superior capability 922 in preserving the input content, which is known to be the most challenging constraint in Text Style 923 Transfer. In compliance with [13], we also report the Krippendorff's alpha inter-rater agreement 924 coefficient, which equals $\alpha = 0.76$. This high score indicates a significant level of agreement among 925 raters, reinforcing the consistency of the conducted human evaluation. 926

Results on Mixed-Style Inputs 5.9

We evaluated the models' ability to preserve content while transferring style on datasets with inputs composed of mixed-style text segments. The mixed-style text versions are generated by

19

930 931

927

928

929

proportionally appending pieces of text of different styles. Table 10 shows the results on the 932 GYAFC-family dataset with a mix of formal/informal text segments with varying mixing ratios. 933 934 Hereafter we will focus on formality transfer because it is more likely to have mixed-style text than in sentiment transfer cases. Additional results are available in Appendix. 935

Overall, the proposed model achieves the best balance of style accuracy and content preservation 936 across different mixing ratios. On GYAFC-family it obtains the highest geometric and harmonic 937 mean in 5 out of 6 configurations. Notably, the performance is relatively strong even in more mixed 938 939 settings such as 50% - 50%, demonstrating its ability to effectively disentangle and transfer style at the segment level rather than just averaging effects across the full input. 940

The baseline models, namely CrossAlignment and MultiDecoder, exhibited consistently lower 941 performance, leading to a notable decrease in overall effectiveness in mixed settings. The latent space variant of CycleGAN also lagged behind our approach, highlighting the benefit of applying 944 the cycle-consistency constraint directly to the raw input sequences. Llama consistently came second to our proposed approach. However, its performance degraded more significantly than 945 CycleGAN as the mixture became more balanced. This suggests CycleGAN may have an advantage 946 in more ambiguous scenarios where the overall style is unclear. 947

These results demonstrate that the proposed methodology is highly effective at style transfer even when the input text contains mixtures of different styles, outperforming prior work especially on more balanced mixtures. This underscores its ability to perform style transfer at a fine-grained segment level.

5.10 Ablation studies

In this section, we delve into the results of two complementary ablation studies. The first experiment explores the impact of the λ_{cuc} scaling factor in the cycle-consistency loss, whereas the second one analyzes the effect of the pre-trained style classifier, considering both the additional loss term and the model used.

Cycle-consistency loss. In this ablation study, we investigate the impact on performance when varying the cycle-consistency loss coefficient λ_{cyc} . To better analyze and isolate the effect of this loss component, we conduct this analysis by excluding the pre-trained style classifier and its corresponding loss term. Consequently, we set $\lambda_{gen} = \lambda_{dis} = 1$, $\lambda_{style} = 0$ and experiment with different values of $\lambda_{cuc} \in \{0, 0.1, 1, 10, 50, 100\}$. Additional results for the BART-base model on the GYAFC-family dataset can be found in Appendix.

In general, the values of the ref-BLEU and style accuracy metrics increase while increasing the value of λ_{cyc} . However, the ref-BLEU increase appears to be quite limited for large λ_{cyc} values, while accuracy still exhibits some room for improvement. Notably, disabling the cycle-consistency loss (i.e., setting $\lambda_{cyc} = 0$) results in a significant performance drop in terms of ref-BLEU (i.e., -4.2 compared to $\lambda_{cuc} = 0.1$). The performance drop is even more pronounced in terms of style accuracy (-6.1). These results confirm the importance of the cycle-consistency loss.

Pre-trained style classifier. In this ablation study, we aim to analyze the impact of the pre-trained 971 style classifier. By enabling/disabling the classifier component, we introduce or eliminate the 972 classifier loss contribution (see Equations 4 and 5). To ensure a complete overview of the classifier 973 contribution, we averaged the evaluation metrics reported in Figure 3 across all the models trained 974 on each dataset. The results reported in Figure 3 show that the introduction of the pre-trained 975 classifier in the training process has a positive impact on all four evaluation metrics. In terms of 976 BLEU scores, it achieves negligible improvements. However, the style classifier yields a +1.0 BLEU 977 score improvement in the music domain of the GYAFC dataset. The limited effect on BLEU can 978 be motivated by the fact that the pre-trained style classifier's objective is to improve the style 979

948

949

950

951 952

953 954

955

956

957 958

959

960

961

962

963

964

965

966

967

968

969

970

980

20

ACM Trans. Intell. Syst. Technol., Vol. 1, No. 1, Article . Publication date: July 2024.

981	Table 10. Results on the GYAFC-family dataset with mixed style for different mixing ratios - ref-BLEU (ref-B
982	avg), style accuracy with BERT and TextCNN (acc _{BERT} , acc _{CNN}), geometric mean (GM) and harmonic mean
983	(HM) of ref-BLEU and style accuracy.

984						~ ~ ~		1	~ ~ ~	
985		mixing	model	ref-B avg	acc _{BERT}	GM	HM	acc _{CNN}	GM	HM
986			CrossAlignment [34]	9.1	50.9	21.5	15.4	74.9	26.1	16.2
087		0.400	MultiDecoder [5]	18.0	50.8	30.2	26.6	69.2	35.3	28.6
967		0-100	CycleGAN latent 15 (large)	15.6	44.5	26.3	23.1	37.4	24.2	22.0
988			CycleGAN 15 (large)	97.7	95.5	96.6	96.6	96.3	97.0	97.0
989	-		Liama2-/B-Chat	92.0	96.5	94.2	94.2	96.9	94.4	94.4
990			CrossAlignment [54]	3.2	07.0	14.0	0.1 5.0	84.0	10.4	0.2
991		05 75	MultiDecoder [5]	3.2	42.0	11.6	5.9	46.0	12.1	6.0
002		25-75	CycleGAN latent 15 (large)	14.2	51.9	27.1	22.3	52.7	27.4	22.4 75.4
<i>772</i>			CycleGAN 15 (large)	62.Z	98.3	78.2	76.2	95.8	77.2	75.4
993	-		Liama2-/B-Chat	58.4	97.0	/5.3	72.9	92.6	/3.5	/1.0
994			CrossAlignment [34]	5.8	67.2	16.0	7.2	82.7	1/./	7.3
995		00.44	MultiDecoder [5]	4.0	46.1	13.6	7.4	48.7	14.0	7.4
996		33-66	CycleGAN latent 15 (large)	15.5	49.1	27.6	23.6	49.9	27.8	23.7
997			CycleGAN 15 (large)	65.6	97.3	79.9	78.3	94.3	78.7	77.4
			Liama2-/B-Chat	58.1	97.3	/5.2	/2.8	93.9	/3.9	/1.8
998			CrossAlignment [34]	4.5	61.1	10.0	8.4	//.0	18.7	8.5
999		50 50	MultiDecoder [5]	4.0	55./	12.5	8.1	53.4	12.4	8.1
1000		50-50	CycleGAN latent 15 (large)	15.8	02.2	29.5	24.5	52.7	28.9	24.5
1001			Lieme 2 7B Chet	00.7	92.5	79.0	70.0	91.0	79.1	70.3
1002			Crease Alignment [24]	05./	90.4	/5.9	/4./	81.0	/4.0	/3.2
1002			CrossAlignment [54]	3.5	02.2	14.8	0.0	81.9 42.5	10.9	0.7
1005		66 22	CycleCAN latent TE (large)	3.0 15 5	42.5	12.4	0.0	45.5	12.5	0.0
1004		00-33	CycleGAN latent 15 (large)	15.5	45.9	20.7	23.2	45.9	20.7	23.2
1005			Lloma 2 7B Chat	66.0	90.0 01.0	78.0	77.0	91.7	75.5	74.0
1006			Cross Alignment [24]	00.9	51.0	12.0	- / /.1 	03.2	15.2	-/4.7 5.4
1007			MultiDagadan [5]	2.0	40.7	15.0	5.4	03.0	13.5	5.4
1008		75 25	CycleCAN latent T5 (large)	14.1	42.7	25.0	0.4 01 4	42.7	25.6	0.4 01.6
1000		75-25	CycleGAN latent 15 (large)	64.0	44.5	25.0 75 5	21.4	40.4	25.0	21.0
1009			Llama 2 7B Chat	63.5	00.7	73.3	76.9	95.4 86.0	743	73.4
1010			Cross Alignment [24]	7 5	93.2	26.6	12.0	84.0	74.5	12.9
1011			MultiDecoder [5]	/.J	94.J 88 J	20.0	13.9	74.1	23.1	13.0
1012		100-0	CycleGAN latent T5 (large)	15.1	00.2 01 /	34.0	22.0 25.0	/4.1 883	36.5	22.3 25.8
1013		100-0	CycleGAN T5 (large)	02.1	91.4	57.2 04 F	43.9 04 E	00.3	30.3 02.9	23.0 03.9
1014			Llama2-7B-Chat	017	97.0	94.3 02.5	94.) 02.5	95.5	9 3.0 01.2	9 3.0 01.2
1014				91./	93.3	94.3	94.3	90.7	91.2	91.4
1015	60		80		70			70		



Fig. 3. Effect of the pre-trained style classifier on the evaluation metrics across different datasets. Results are averaged over all the tested models.

transfer accuracy, and thus, it does not necessarily affect the BLEU scores. On the contrary, we 1030 observe remarkable improvements in terms of style accuracy on the two domains of the GYAFC 1031 1032 dataset. Specifically, we achieve an absolute gain of +4.0 and +7.9 points in accuracy scores, which corresponds to the relative improvements of +8.8% and +18.4% when compared to the classifier-free 1033 counterparts. Finally, by analyzing the impact on the geometric mean and harmonic mean of 1034 BLEU and style accuracy, we can observe an overall improvement of up to +4.1 and +3.7 points, 1035 respectively. The geometric mean and harmonic mean provide a more comprehensive evaluation 1036 of the overall performance of the approach, taking into account the trade-off between the two 1037 separate metrics. These results, therefore, confirm the effectiveness of the pre-trained classifier in 1038 enhancing the quality of the generated text. 1039

The evaluation results show a surprising lack of performance improvement on the Yelp dataset. 1040 One possible explanation for this phenomenon is that the sentiment transfer task already achieves 1041 high style accuracy scores, even without the pre-trained classifier. This may suggest that the model's 1042 pre-existing capability to perform style transfer is already sufficient to achieve high accuracy scores, 1043 making the pre-trained classifier's contribution negligible in this case. Also, the larger size of the 1044 Yelp dataset may already provide the model with a sufficient amount of training data to effectively 1045 capture style transfer patterns. Moreover, as described in Section 5.3, since the Yelp dataset does not 1046 include human references for the validation set, style accuracy is already taken into account when 1047 performing the hyperparameter tuning and selecting the best checkpoints. This may be another 1048 possible explanation for the limited impact of the pre-trained classifier on this dataset. 1049

As the quality of the pre-trained style classifier may affect the overall performance of our proposed architecture, we extend this ablation study by also testing other style classifiers. In addition to the BERT-base model used in our architecture, we test the following models: BERT-large, RoBERTa-base [20], RoBERTa-large, and DistilBERT-base. More detailed results for the BART-base model on the GYAFC-family dataset can be found in Appendix.

Overall, we observe that the specific style classifier chosen does not have a significant impact on 1055 performance. Specifically, the differences among the various classifiers in ref-BLEU are negligible 1056 (i.e., ± 0.1), and similarly for style accuracy, where fluctuations range from ± 0.4 to ± 2.5 . The largest 1057 differences in performance are observed with the DistilBERT-base model, showing a drop of -1.0 1058 and -4.6 in ref-BLEU and accuracy scores, respectively. This result is expected, given that the 1059 DistilBERT-base model is the lightest among those tested. Nevertheless, all tested models perform 1060 generally well, indicating that the quality of the pre-trained style classifier has a limited impact on 1061 the final performance. 1062

1064 6 CONCLUSIONS AND FUTURE WORK

1063

1078

In this paper, we presented a new approach to self-supervised Text Style Transfer using Cycleconsistent Generative Adversarial Networks (CycleGANs). Thanks to the joint use of cycle-consistency and a pre-trained style classification loss, our method is able to effectively transfer the style of a source text to a target text without the need for labeled data. The experimental results, achieved on three benchmark datasets and two different TST tasks, show that our method performs better than state-of-the-art approaches in terms of quality of the generated text and ability to preserve the content of the source text, particularly when mixed-style inputs are processed.

Limitations. The application of the proposed approach to real case studies should consider the following potential limitations. (1) We made the assumption that the source and target domains have approximately the same distributions. As a result, the model could be unable to learn the correct mapping between the two style attributes if this assumption does not hold true. (2) The presented approach may be misused to maliciously manipulate the text style or sentiment. For example, by transferring the style of credible sources to untruthful content, the proposed method

ACM Trans. Intell. Syst. Technol., Vol. 1, No. 1, Article . Publication date: July 2024.

might be employed to automatically generate fake news or propaganda. (3) The currently proposed
architecture handles one specific pair of source and target styles. This leads to scalability issues, as
it would require training a separate architecture for each new pair of styles.

Despite the aforementioned limitations, the usability of the proposed method is quite promising in various real-world scenarios. With a responsible deployment and a careful consideration of the main ethical concerns, our approach can relevantly contribute to the advancement of the TST research field and enable innovative NLP applications in fields such as marketing, content generation, and digital storytelling.

Future work. We plan to extend our work across multiple directions. (1) We aim to expand the 1087 capabilities of the proposed architecture by investigating its performance in a multilingual setting. 1088 1089 Transferring the style attributes across languages is potentially challenging as entails not only 1090 capturing stylistic nuances but also handling language-specific characteristics. By considering this 1091 aspect, we can evaluate the model's ability to generalize and adapt to diverse linguistic contexts. (2) The flexibility of our method allows us to explore its applicability to new domains and tasks. 1092 For instance, we would like to further explore the following two related tasks: Aspect-level style 1093 transfer [13], and Controllable text generation [10]. In both cases, the goal is to selectively transfer 1094 specific attributes or aspects of the writing style while preserving the rest. (3) We also envisage to 1095 extend our approach to handle more than a single pair of styles simultaneously (see Section 4.5). 1096 1097

1098 ACKNOWLEDGMENTS

1099 The work by Giuseppe Gallipoli was carried out within the MICS (Made in Italy – Circular and Sus-1100 tainable) Extended Partnership and received funding from the European Union Next-GenerationEU 1101 (PIANO NAZIONALE DI RIPRESA E RESILIENZA (PNRR) - MISSIONE 4 COMPONENTE 2, IN-1102 VESTIMENTO 1.3 - D.D. 1551.11-10-2022, PE00000004). This study was also partially carried out 1103 within the FAIR (Future Artificial Intelligence Research) and received funding from the European 1104 Union Next-GenerationEU (PIANO NAZIONALE DI RIPRESA E RESILIENZA (PNRR) - MISSIONE 1105 4 COMPONENTE 2, INVESTIMENTO 1.3 - D.D. 1555.11-10-2022, PE00000013). This manuscript 1106 reflects only the authors' views and opinions, neither the European Union nor the European 1107 Commission can be considered responsible for them.

1109 **REFERENCES**

- Liqun Chen et al. 2018. Adversarial Text Generation via Feature-Mover's Distance. In Advances in Neural Information Processing Systems 31: Annual Conference on Neural Information Processing Systems 2018, NeurIPS 2018, December 3-8, 2018, Montréal, Canada, Vol. 31. Curran Associates, Inc., 4671–4682. https://proceedings.neurips.cc/paper/2018/hash/ 074177d3eb6371e32c16c55a3b8f706b-Abstract.html
- [2] Hyung Won Chung et al. 2024. Scaling instruction-finetuned language models. *Journal of Machine Learning Research* 25, 70 (2024), 1–53.
- [3] Ning Dai, Jianze Liang, Xipeng Qiu, and Xuanjing Huang. 2019. Style Transformer: Unpaired Text Style Transfer without
 Disentangled Latent Representation. In *Proceedings of the 57th Annual Meeting of the Association for Computational* Linguistics. Association for Computational Linguistics, Florence, Italy, 5997–6007. https://doi.org/10.18653/v1/P19-1601
- [4] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. In Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers). Association for Computational Linguistics, Minneapolis, Minnesota, 4171–4186. https://doi.org/10.18653/v1/N19-1423
- 1121
 [5] Zhenxin Fu et al. 2018. Style Transfer in Text: Exploration and Evaluation. Proceedings of the AAAI Conference on Artificial Intelligence 32, 1 (Apr. 2018). https://doi.org/10.1609/aaai.v32i1.11330
- [6] Ian Goodfellow et al. 2014. Generative Adversarial Nets. In Advances in Neural Information Processing Systems, Vol. 27.
 Curran Associates, Inc., Montréal. https://proceedings.neurips.cc/paper/2014/file/5ca3e9b122f61f8f06494c97b1afccf3-Paper.pdf
- 1125[7] Junxian He, Xinyi Wang, Graham Neubig, and Taylor Berg-Kirkpatrick. 2020. A Probabilistic Formulation of Unsu-
pervised Text Style Transfer. In 8th International Conference on Learning Representations, ICLR 2020. Addis Abeba.
- 1127

1128 https://arxiv.org/abs/2002.03912

- [8] Zhiting Hu et al. 2018. Toward Controlled Generation of Text. arXiv:1703.00955 [cs.LG]
- [9] Zhiqiang Hu, Roy Ka-Wei Lee, Charu C. Aggarwal, and Aston Zhang. 2022. Text Style Transfer: A Review and Experimental Evaluation. *SIGKDD Explor*. 24, 1 (2022), 14–45. https://doi.org/10.1145/3544903.3544906
- [10] Zhiting Hu and Li Erran Li. 2021. A Causal Lens for Controllable Text Generation. In Advances in Neural Information
 Processing Systems 34: Annual Conference on Neural Information Processing Systems 2021, NeurIPS 2021, December
 6-14, 2021, virtual. 24941–24955. https://proceedings.neurips.cc/paper/2021/hash/d0f5edad9ac19abed9e235c0fe0aa59f Abstract.html
- [11] Fei Huang et al. 2021. NAST: A Non-Autoregressive Generator with Word Alignment for Unsupervised Text Style Transfer. In *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*. Association for Computational Linguistics, Online, 1577–1590. https://doi.org/10.18653/v1/2021.findings-acl.138
- [12] Yufang Huang et al. 2020. Cycle-Consistent Adversarial Autoencoders for Unsupervised Text Style Transfer. In *Proceedings of the 28th International Conference on Computational Linguistics*. International Commutational Linguistics, Barcelona, Spain (Online), 2213–2223. https://doi.org/10.18653/v1/2020.coling-main.201
- [13] Di Jin et al. 2022. Deep Learning for Text Style Transfer: A Survey. *Computational Linguistics* 48, 1 (March 2022), 155–205. https://doi.org/10.1162/coli_a_00426
 [141] Di Jin et al. 2022. Deep Learning for Text Style Transfer: A Survey. *Computational Linguistics* 48, 1 (March 2022), 155–205. https://doi.org/10.1162/coli_a_00426
- [14] Nitish Shirish Keskar et al. 2019. CTRL: A Conditional Transformer Language Model for Controllable Generation.
 CoRR abs/1909.05858 (2019). http://arxiv.org/abs/1909.05858
- [15] Yoon Kim. 2014. Convolutional Neural Networks for Sentence Classification. In *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*. Association for Computational Linguistics, Doha, Qatar, 1746–1751. https://doi.org/10.3115/v1/D14-1181
- [16] Guillaume Lample et al. 2019. Multiple-Attribute Text Rewriting. In 7th International Conference on Learning Representations, ICLR 2019, New Orleans, LA, USA, May 6-9, 2019. OpenReview.net. https://openreview.net/forum?id=H1g2NhC5KQ
- [17] Mike Lewis et al. 2020. BART: Denoising Sequence-to-Sequence Pre-training for Natural Language Generation, Translation, and Comprehension. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*. Association for Computational Linguistics. Online, 7871–7880. https://doi.org/10.18653/v1/2020.aclmain.703
- [18] Juncen Li, Robin Jia, He He, and Percy Liang. 2018. Delete, Retrieve, Generate: a Simple Approach to Sentiment and
 Style Transfer. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*. Association for Computational Linguistics, New
 Orleans, Louisiana, 1865–1874. https://doi.org/10.18653/v1/N18-1169
- [19] Dayiheng Liu et al. 2019. Revision in Continuous Space: Fine-Grained Control of Text Style Transfer. CoRR abs/1905.12304 (05 2019). https://arxiv.org/abs/1905.12304
- [20] Yinhan Liu et al. 2019. RoBERTa: A Robustly Optimized BERT Pretraining Approach. *arXiv preprint arXiv:1907.11692* (2019).
- 1158[21] Michela Lorandi et al. 2023. Adapting the CycleGAN architecture for text style transfer. (2023). https://doi.org/10.11595281/zenodo.8268839
- [22] Ilya Loshchilov and Frank Hutter. 2019. Decoupled Weight Decay Regularization. In 7th International Conference on Learning Representations, ICLR 2019, New Orleans, LA, USA, May 6-9, 2019. OpenReview.net. https://openreview.net/ forum?id=Bkg6RiCqY7
- [23] Fuli Luo et al. 2019. A Dual Reinforcement Learning Framework for Unsupervised Text Style Transfer. In *Proceedings* of the 28th International Joint Conference on Artificial Intelligence, IJCAI 2019. Macao. https://arxiv.org/abs/1905.10060
- [24] Fatemehsadat Mireshghallah, Kartik Goyal, and Taylor Berg-Kirkpatrick. 2022. Mix and Match: Learning-free Controllable Text Generationusing Energy Language Models. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. Association for Computational Linguistics, Dublin, Ireland, 401–415. https://doi.org/10.18653/v1/2022.acl-long.31
- [25] Matt Post. 2018. A Call for Clarity in Reporting BLEU Scores. In *Proceedings of the Third Conference on Machine Translation: Research Papers*. Association for Computational Linguistics, Brussels, Belgium, 186–191. https://doi.org/10.18653/v1/W18-6319
- [26] Shrimai Prabhumoye, Yulia Tsvetkov, Ruslan Salakhutdinov, and Alan W Black. 2018. Style Transfer Through Back-Translation. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. Association for Computational Linguistics, Melbourne, Australia, 866–876. https://doi.org/10.18653/v1/P18-1080
- [173 [27] Colin Raffel et al. 2020. Exploring the Limits of Transfer Learning with a Unified Text-to-Text Transformer. *Journal of Machine Learning Research* 21, 140 (2020), 1–67. http://jmlr.org/papers/v21/20-074.html
- 1175 1176
- ACM Trans. Intell. Syst. Technol., Vol. 1, No. 1, Article . Publication date: July 2024.

- [1177 [28] Sudha Rao and Joel Tetreault. 2018. Dear Sir or Madam, May I Introduce the GYAFC Dataset: Corpus, Benchmarks
 and Metrics for Formality Style Transfer. In *Proceedings of the 2018 Conference of the North American Chapter of the* Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers). Association for
 Computational Linguistics, New Orleans, Louisiana, 129–140. https://doi.org/10.18653/v1/N18-1012
- [29] Machel Reid and Victor Zhong. 2021. LEWIS: Levenshtein Editing for Unsupervised Text Style Transfer. In *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021.* Association for Computational Linguistics, Online, 3932–3944. https://doi.org/10.18653/v1/2021.findings-acl.344
- [30] Emily Reif et al. 2022. A Recipe for Arbitrary Text Style Transfer with Large Language Models. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*. Association for Computational Linguistics, Dublin, Ireland, 837–848. https://doi.org/10.18653/v1/2022.acl-short.94
- [31] Abhilasha Sancheti, Kundan Krishna, Balaji Vasan Srinivasan, and Anandhavelu Natarajan. 2020. Reinforced Rewards
 Framework for Text Style Transfer. In Advances in Information Retrieval: 42nd European Conference on IR Research,
 ECIR 2020, Proceedings, Part I. Lisbon. https://arxiv.org/abs/2005.05256
- 1188[32]Victor Sanh, Lysandre Debut, Julien Chaumond, and Thomas Wolf. 2019. DistilBERT, a distilled version of BERT:
smaller, faster, cheaper and lighter. *arXiv preprint arXiv:1910.01108* (2019).
- [33] Mingyue Shang et al. 2019. Semi-supervised Text Style Transfer: Cross Projection in Latent Space. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*. Association for Computational Linguistics, Hong Kong, China, 4937–4946. https://doi.org/10.18653/v1/D19-1499
- [34] Tianxiao Shen, Tao Lei, Regina Barzilay, and Tommi S. Jaakkola. 2017. Style Transfer from Non-Parallel Text by Cross-Alignment. In *Proceedings of the 31st International Conference on Neural Information Processing Systems*. Long Beach, California. https://arxiv.org/abs/1705.09655
- [35] Yukai Shi et al. 2021. GTAE: Graph Transformer–Based Auto-Encoders for Linguistic-Constrained Text Style Transfer.
 ACM Trans. Intell. Syst. Technol. 12, 3, Article 32 (jun 2021), 16 pages. https://doi.org/10.1145/3448733
- [197 [36] Martina Toshevska and Sonja Gievska. 2022. A Review of Text Style Transfer Using Deep Learning. *IEEE Transactions* on Artificial Intelligence 3, 5 (2022), 669–684. https://doi.org/10.1109/TAI.2021.3115992
- [37] Hugo Touvron et al. 2023. Llama 2: Open foundation and fine-tuned chat models. arXiv preprint arXiv:2307.09288 (2023).
 [200 [2023].
- [38] Ashish Vaswani et al. 2017. Attention is All you Need. In Advances in Neural Information Processing Systems, Vol. 30.
 [201] Curran Associates, Inc. https://proceedings.neurips.cc/paper/2017/file/3f5ee243547dee91fbd053c1c4a845aa-Paper.pdf
- 1202[39]Yunli Wang et al. 2019. Harnessing Pre-Trained Neural Networks with Rules for Formality Style Transfer. In Proceedings1203of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference1204on Natural Language Processing (EMNLP-IJCNLP). Association for Computational Linguistics, Hong Kong, China,12053573–3578. https://doi.org/10.18653/v1/D19-1365
- [40] Thomas Wolf et al. 2020. Transformers: State-of-the-Art Natural Language Processing. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*. Association for Computational Linguistics, Online, 38–45. https://doi.org/10.18653/v1/2020.emnlp-demos.6
- [41] Yujia Xie, Xiangfeng Wang, Ruijia Wang, and Hongyuan Zha. 2019. A Fast Proximal Point Method for Computing Exact Wasserstein Distance. In Proceedings of the Thirty-Fifth Conference on Uncertainty in Artificial Intelligence, UAI 2019, Tel Aviv, Israel, July 22-25, 2019 (Proceedings of Machine Learning Research, Vol. 115). AUAI Press, 433–453. http://proceedings.mlr.press/v115/xie20b.html
- [42] Jingjing Xu et al. 2018. Unpaired Sentiment-to-Sentiment Translation: A Cycled Reinforcement Learning Approach.
 In Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers).
 Association for Computational Linguistics, Melbourne, Australia, 979–988. https://doi.org/10.18653/v1/P18-1090
- [43] Wei Xu et al. 2012. Paraphrasing for Style. In *Proceedings of COLING 2012*. The COLING 2012 Organizing Committee, Mumbai, India, 2899–2914. https://aclanthology.org/C12-1177
 [43] Wei Xu et al. 2012. Paraphrasing for Style. In *Proceedings of COLING 2012*. The COLING 2012 Organizing Committee, Mumbai, India, 2899–2914. https://aclanthology.org/C12-1177
- [44] Yi Zhang, Jingjing Xu, Pengcheng Yang, and Xu Sun. 2018. Learning Sentiment Memories for Sentiment Modification
 without Parallel Data. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*.
 Association for Computational Linguistics, Brussels, Belgium, 1103–1108. https://doi.org/10.18653/v1/D18-1138
- 1218[45]Zhirui Zhang et al. 2018. Style Transfer as Unsupervised Machine Translation. CoRR abs/1808.07894 (08 2018).1219https://arxiv.org/abs/1808.07894
- [46] Yanpeng Zhao et al. 2018. Language Style Transfer from Sentences with Arbitrary Unknown Styles. CoRR abs/1808.04071 (2018). http://arxiv.org/abs/1808.04071
- [47] Jun-Yan Zhu, Taesung Park, Phillip Isola, and Alexei A Efros. 2017. Unpaired Image-to-Image Translation using
 Cycle-Consistent Adversarial Networks. In Computer Vision (ICCV), 2017 IEEE International Conference on. IEEE
 Computer Society, Venice, 2242–2251. https://doi.org/10.1109/ICCV.2017.244
- 1224 1225

ACM Trans. Intell. Syst. Technol., Vol. 1, No. 1, Article . Publication date: July 2024.

M. La Quatra, G. Gallipoli, and L. Cagliero

1226 APPENDIX

- ¹²²⁷ This document contains the following appendices:
- A) Hyperparameter settings;
- B) Additional results on formality transfer;
- C) Additional results on mixed-style inputs;
- D) Additional results on cycle-consistency loss;
- E) Additional results on classifier-guided loss;
- F) Qualitative examples.

1236 A HYPERPARAMETER SETTINGS

In Table 11 we report the optimal hyperparameter values used throughout the experiments on boththe GYAFC and Yelp datasets.

Table 11. Optimal hyperparameter configurations for each dataset and model used in experiments.

	dataset	generator model	learning rate	λ_{cyc}
		BART-base	$5 \cdot 10^{-5}$	
		BART-large	$5 \cdot 10^{-5}$	
	GYAFC-family	T5-small	$5 \cdot 10^{-5}$	10
		T5-base	$5 \cdot 10^{-5}$	
		T5-large	$5 \cdot 10^{-5}$	
		BART-base	$5 \cdot 10^{-5}$	
	GYAFC-music	BART-large	$1 \cdot 10^{-5}$	
		T5-small	$5 \cdot 10^{-5}$	1
		T5-base	$5 \cdot 10^{-5}$	
		T5-large	$5 \cdot 10^{-5}$	
		BART-base	$1\cdot 10^{-4}$	
		BART-large	$1 \cdot 10^{-5}$	
	Yelp	T5-small	$1 \cdot 10^{-3}$	10
	-	T5-base	$1 \cdot 10^{-4}$	
		T5-large	$5 \cdot 10^{-5}$	

1258 1259

1260 1261

1262

1263

1264

1265

1266

B ADDITIONAL RESULTS ON FORMALITY TRANSFER

In this section, we report the additional results for the formality transfer task. Specifically, Table 12 shows the detailed results on the GYAFC-family dataset in the informal-to-formal style transfer direction. Even when we restrict the analysis to a specific style transfer direction, our proposed approach achieves the best performance (i.e., +10.3 ref-BLEU).

Table 13 shows the detailed results on the GYAFC-music dataset in the informal-to-formal style transfer direction. Similar to the GYAFC-family domain, our method largely outperforms the other approaches (i.e., +34.1 ref-BLEU).

Table 14 includes the detailed results obtained by Llama2 [37] on the GYAFC-family and GYAFCmusic datasets in the informal-to-formal style transfer direction. The best ref-BLEU performance is achieved for k = 5 and k = 3, respectively. However, our approach confirms its superior performance by a large margin (i.e., +18.6 and +13.1 ref-BLEU, respectively).

1274

26

1235

1239 1240

1270		ref-B	ACCREPT	GM	НМ	acconn	GM	НМ
1279	RetrieveOnly [18]	13	92 5	10.0	8.2	7.5	5.7	5.5
1280		4.5	92.5	17.7	5.2	7.5	5.7	5.5
1281	Back Translation [26]	5.5	6.4	5.9	5.9	8.3	6.8	6.6
1282	StyleEmbedding [5]	12.0	10.8	11.4	11.4	0.0	0.0	0.0
1283	MultiDecoder [5]	17.3	5.0	9.3	7.8	0.0	0.0	0.0
1284	CrossAlignment [34]	8.1	45.1	19.1	13.7	5.7	6.8	6.7
1285	UnpairedTranslation [42]	4.3	36.9	12.6	7.7	7.5	5.7	5.5
1286	DeleteOnly [18]	35.1	6.1	14.6	10.4	0.2	2.6	0.4
1287	DeleteAndRetrieve [18]	24.4	33.6	28.6	28.3	4.6	10.6	7.7
1288	TemplateBased [18]	39.9	16.3	25.5	23.1	5.2	14.4	9.2
1289	UnsupervisedMT [45]	37.9	59.3	47.4	46.2	3.8	12.0	6.9
1290	DualRL [23]	51.6	28.7	38.5	36.9	0.8	6.4	1.6
1291	NASTLatentLearn [11]	51.1	37.6	43.8	43.3	30.0	39.1	37.8
1292	CycleGAN BART (base)	58.7	42.2	49.8	49.1	47.0	52.5	52.2
1293	CycleGAN BART (large)	59.3	41.2	49.4	48.6	46.2	52.3	51.9
1294	CycleGAN T5 (small)	54.3	28.7	39.5	37.6	34.8	43.5	42.4
1295	CycleGAN T5 (base)	56.7	28.9	40.5	38.3	32.1	42.7	41.0
1296	CycleGAN T5 (large)	61.9	49.2	55.2	54.8	53.2	57.4	57.2

Table 12. Results on the GYAFC-family dataset | informal \rightarrow formal - ref-BLEU (ref-B), style accuracy with BERT and TextCNN (acc_{BERT}, acc_{CNN}), geometric mean (GM) and harmonic mean (HM) of ref-BLEU and style accuracy.

Table 13. Results on the GYAFC-music dataset | informal \rightarrow formal – ref-BLEU (ref-B), style accuracy with BERT and TextCNN (acc_{BERT}, acc_{CNN}), geometric mean (GM) and harmonic mean (HM) of ref-BLEU and style accuracy. * denotes results from the paper.

	ref-B	acc_{BERT}	GM	HM	acc _{CNN}	GM	HM
DeepLatent [7]	26.4	77.8	45.3	39.4	25.6	26.0	26.0
ReinfRewards* [31]	28.6	72.3	45.4	41.0	-	-	-
MixAndMatch* [24]	-	19.0	-	-	-	-	-
CycleGAN BART (base)	57.6	53.7	55.6	55.6	48.4	52.8	52.6
CycleGAN BART (large)	55.3	42.8	48.7	48.3	38.4	46.1	45.3
CycleGAN T5 (small)	51.6	34.4	42.1	41.3	29.5	39.0	37.5
CycleGAN T5 (base)	55.2	36.5	44.9	43.9	33.4	42.9	41.6
CycleGAN T5 (large)	62.7	67.1	64.9	64.8	61.6	62.1	62.1

C ADDITIONAL RESULTS ON MIXED-STYLE INPUTS

1316Table 15 reports the results for the GYAFC-music formality transfer dataset in the mixed-style1317scenario. Our approach demonstrates superior performance to the other competitors in terms of ref-1318BLEU, geometric and harmonic means across all mixing ratios. Notably, the ref-BLEU performance1319gap with the second-best performer (i.e., Llama2) ranges from +8.4 in the 25 - 75 case to +22.2 in1320the 50 - 50 case. These results, similar to those observed in the GYAFC-family dataset, underscore1321the enhanced capability of our method in handling mixed-style texts, especially those without a1322predominant style.

Table 14. Results on the GYAFC-family and GYAFC-music datasets | informal \rightarrow formal of Llama2-7B-Chat model for varying number of examples *k* in the 0/few-shot setting – ref-BLEU (ref-B avg), style accuracy with BERT and TextCNN (acc_{BERT}, acc_{CNN}), geometric mean (GM) and harmonic mean (HM) of ref-BLEU and style accuracy.

dataset	k	ref-B	acc_{BERT}	GM	HM	acc _{CNN}	GM	HM
	0	23.1	81.8	43.4	36.0	84.2	44.1	36.2
	1	20.9	90.8	43.5	34.0	92.2	43.9	34.1
GYAFC-family	3	41.7	90.8	61.5	57.1	91.3	61.7	57.2
	5	43.3	89.5	62.2	58.3	90.1	62.4	58.5
	10	35.9	92.6	57.6	51.7	93.1	57.8	51.8
	0	27.5	87.4	49.0	41.8	80.4	47.0	41.0
	1	39.2	92.2	60.1	55.0	87.4	58.5	54.1
GYAFC-music	3	49.6	91.0	67.2	64.2	87.0	65.7	63.2
	5	44.7	91.9	64.1	60.1	87.9	62.7	59.2
	10	46.2	93.0	65.5	61.7	88.9	64.1	60.8

D ADDITIONAL RESULTS ON CYCLE-CONSISTENCY LOSS

Table 16 reports the detailed results of the ablation study on the effect of the λ_{cyc} hyperparameter. Notably, the best performance in terms of ref-BLEU and style accuracy is achieved for higher λ_{cyc} values. This underscores the importance of the cycle-consistency loss and emphasizes the need for accurate tuning of the corresponding scaling factor.

E ADDITIONAL RESULTS ON CLASSIFIER-GUIDED LOSS

In this section, we report in Tables 17-21 the detailed set of results of the ablation study analyzing
the effect of the pre-trained style classifier. Additionally, the results of the ablation study on the
model used as pre-trained style classifier are presented in Table 22.

More specifically, Tables 17 and 18 report the results achieved on the GYAFC-family dataset. 1352 Almost all the tested models benefit from the introduction of the pre-trained classifier. Notably, 1353 the best-performing model (T5 large) achieves an improvement up to +1.0 and +10.3 points in 1354 the BLEU score and style accuracy, respectively. The accuracy improvement is even higher while 1355 considering only the informal-to-formal direction (i.e., +16.9 points). Similar observations hold for 1356 the results reported in Tables 19 and 20 (e.g., +27.5 points in accuracy) which display the results 1357 on the music domain. Finally, Table 21 shows the results on the Yelp dataset. Although the best 1358 BLEU and accuracy scores are achieved by the models trained without the style classifier, it can 1359 be noticed that the introduction of the classifier-guided loss yields an improvement in the overall 1360 performance score, as indicated by the geometric mean and harmonic mean. 1361

F QUALITATIVE EXAMPLES

1364 In this section, we show qualitative examples of pairs of original and transformed texts.

To provide a more comprehensive evaluation of the quality of the generated text, we conduct a qualitative analysis by comparing it with both the ground truth and a subset of the baseline methods. This analysis allows us to identify the success and failure cases of our approach more tangibly and compare them with the outputs of other state-of-the-art methods. Tables 23 and 25 report a selection of success and failure cases for the formality transfer task within the family domain, Tables 24 and 26 contain the results in the music domain, whereas Tables 27 and 28 report qualitative results for the sentiment transfer task.

1372

1362

1363

28

2	റ
	ч
~	~

1373	Table 15. Results on the GYAFC-music dataset with mixed style for different mixing ratios - ref-BLEU (ref-B
1374	avg), style accuracy with BERT and TextCNN (acc _{BERT} , acc _{CNN}), geometric mean (GM) and harmonic mean
1375	(HM) of ref-BLEU and style accuracy.

		0	DLICI			CININ		1 11 11
	CrossAlignment [34]	8.6	49.9	20.7	14.7	52.9	21.3	14.8
	MultiDecoder [5]	13.2	65.4	29.4	22.0	71.3	30.7	22.3
0-100	CycleGAN latent T5 (large)	16.8	69.1	34.1	27.0	28.3	21.8	21.1
	CycleGAN T5 (large)	97.9	98.7	98.3	98.3	95.7	96.8	96.8
	Llama2-7B-Chat	80.4	98.9	89.2	88.7	97.1	88.4	88.0
	CrossAlignment [34]	2.7	66.2	13.4	5.2	61.2	12.9	5.2
	MultiDecoder [5]	2.9	54.8	12.6	5.5	55.3	12.7	5.5
25-75	CycleGAN latent T5 (large)	13.9	52.9	27.1	22.0	51.0	26.6	21.8
	CycleGAN T5 (large)	61.1	93.9	75.7	74.0	95.6	76.4	74.6
	Llama2-7B-Chat	52.7	97.1	71.5	68.3	92.0	69.6	67.0
	CrossAlignment [34]	3.2	65.3	14.5	6.1	61.1	14.0	6.1
	MultiDecoder [5]	3.5	51.7	13.5	6.6	51.9	13.5	6.6
33-66	CycleGAN latent T5 (large)	16.3	51.0	28.8	24.7	49.8	28.5	24.6
	CycleGAN T5 (large)	65.5	94.3	78.6	77.3	95.4	7 9.0	77.7
	Llama2-7B-Chat	49.0	93.4	67.7	64.3	89.7	66.3	63.4
	CrossAlignment [34]	3.8	63.4	15.5	7.2	62.8	15.4	7.2
	MultiDecoder [5]	4.1	35.7	12.1	7.4	35.5	12.1	7.4
50-50	CycleGAN latent T5 (large)	16.2	60.2	31.2	25.5	53.6	29.5	24.9
	CycleGAN T5 (large)	69.2	94.3	80.8	79.8	93.0	80.2	79.4
	Llama2-7B-Chat	46.9	89.4	64.8	61.5	86.2	63.6	60.7
	CrossAlignment [34]	3.2	62.9	14.2	6.1	60.5	13.9	6.1
	MultiDecoder [5]	3.2	47.4	12.3	6.0	48.8	12.5	6.0
66-33	CycleGAN latent T5 (large)	15.1	46.8	26.6	22.8	45.9	26.3	22.7
	CycleGAN T5 (large)	69.7	96.7	82.1	81.0	95.2	81.5	80.5
	Llama2-7B-Chat	56.2	96.9	73.8	71.1	94.0	72.7	70.3
	CrossAlignment [34]	2.7	60.9	12.8	5.2	57.8	12.5	5.2
	MultiDecoder [5]	2.5	49.2	11.1	4.8	52.4	11.4	4.8
75-25	CycleGAN latent T5 (large)	13.9	46.9	25.5	21.4	46.3	25.4	21.4
	CycleGAN T5 (large)	68.2	97.5	81.5	80.3	94.6	80.3	79.3
	Llama2-7B-Chat	56.1	94.1	72.7	70.3	91.5	71.6	69.6
	CrossAlignment [34]	6.4	96.7	24.9	12.0	95.8	24.8	12.0
	MultiDecoder [5]	9.1	78.1	26.7	16.3	73.1	25.8	16.2
100-0	CycleGAN latent T5 (large)	15.5	82.9	35.8	26.1	90.8	37.5	26.5
	CycleGAN T5 (large)	92.1	94.8	93.4	93.4	95.8	93.9	93.9
	Llama2-7B-Chat	83.4	89.6	86.4	86.4	90.7	87.0	86.9
	0-100 25-75 33-66 50-50 66-33 75-25 100-0	CrossAlignment [34] MultiDecoder [5] 0-100 CycleGAN latent T5 (large) CycleGAN T5 (large) Llama2-7B-Chat CrossAlignment [34] MultiDecoder [5] 25-75 CycleGAN latent T5 (large) CycleGAN T5 (large) Llama2-7B-Chat CrossAlignment [34] MultiDecoder [5] 33-66 CycleGAN latent T5 (large) CycleGAN T5 (large) Llama2-7B-Chat CrossAlignment [34] MultiDecoder [5] 50-50 CycleGAN latent T5 (large) CycleGAN T5 (large) Llama2-7B-Chat CrossAlignment [34] MultiDecoder [5] 66-33 CycleGAN latent T5 (large) CycleGAN T5 (large) Llama2-7B-Chat CrossAlignment [34] MultiDecoder [5] 66-33 CycleGAN latent T5 (large) CycleGAN T5 (large) Llama2-7B-Chat CrossAlignment [34] MultiDecoder [5] 75-25 CycleGAN latent T5 (large) Llama2-7B-Chat CrossAlignment [34] MultiDecoder [5] 75-25 CycleGAN latent T5 (large) Llama2-7B-Chat CrossAlignment [34] MultiDecoder [5] 100-0 CycleGAN latent T5 (large) Llama2-7B-Chat	$ \begin{array}{c cccc} CrossAlignment [34] & 8.6 \\ MultiDecoder [5] & 13.2 \\ 0-100 & CycleGAN latent T5 (large) & 16.8 \\ CycleGAN T5 (large) & 97.9 \\ Llama2-7B-Chat & 80.4 \\ \hline CrossAlignment [34] & 2.7 \\ MultiDecoder [5] & 2.9 \\ CycleGAN latent T5 (large) & 13.9 \\ CycleGAN T5 (large) & 61.1 \\ Llama2-7B-Chat & 52.7 \\ \hline CrossAlignment [34] & 3.2 \\ MultiDecoder [5] & 3.5 \\ 33-66 & CycleGAN latent T5 (large) & 16.3 \\ CycleGAN T5 (large) & 65.5 \\ Llama2-7B-Chat & 49.0 \\ \hline CrossAlignment [34] & 3.8 \\ MultiDecoder [5] & 4.1 \\ 50-50 & CycleGAN latent T5 (large) & 16.2 \\ CycleGAN T5 (large) & 69.2 \\ Llama2-7B-Chat & 49.0 \\ \hline CrossAlignment [34] & 3.8 \\ MultiDecoder [5] & 4.1 \\ 50-50 & CycleGAN latent T5 (large) & 16.2 \\ CycleGAN T5 (large) & 69.2 \\ Llama2-7B-Chat & 46.9 \\ \hline CrossAlignment [34] & 3.2 \\ MultiDecoder [5] & 3.2 \\ 66-33 & CycleGAN latent T5 (large) & 15.1 \\ CycleGAN T5 (large) & 69.7 \\ Llama2-7B-Chat & 56.2 \\ \hline CrossAlignment [34] & 2.7 \\ MultiDecoder [5] & 2.5 \\ 75-25 & CycleGAN latent T5 (large) & 13.9 \\ CycleGAN T5 (large) & 68.2 \\ Llama2-7B-Chat & 56.1 \\ \hline CrossAlignment [34] & 6.4 \\ MultiDecoder [5] & 9.1 \\ 100-0 & CycleGAN latent T5 (large) & 15.5 \\ CycleGAN T5 (large) & 92.1 \\ Llama2-7B-Chat & 83.4 \\ \hline \end{array}$	$ \begin{array}{c cccc} CrossAlignment [34] & 8.6 & 49.9 \\ MultiDecoder [5] & 13.2 & 65.4 \\ \hline 0-100 & CycleGAN latent T5 (large) & 16.8 & 69.1 \\ CycleGAN T5 (large) & 97.9 & 98.7 \\ Llama2-7B-Chat & 80.4 & 98.9 \\ \hline CrossAlignment [34] & 2.7 & 66.2 \\ MultiDecoder [5] & 2.9 & 54.8 \\ \hline 25-75 & CycleGAN latent T5 (large) & 13.9 & 52.9 \\ CycleGAN T5 (large) & 61.1 & 93.9 \\ Llama2-7B-Chat & 52.7 & 97.1 \\ \hline CrossAlignment [34] & 3.2 & 65.3 \\ MultiDecoder [5] & 3.5 & 51.7 \\ \hline 33-66 & CycleGAN latent T5 (large) & 16.3 & 51.0 \\ CycleGAN T5 (large) & 65.5 & 94.3 \\ Llama2-7B-Chat & 49.0 & 93.4 \\ \hline CrossAlignment [34] & 3.8 & 63.4 \\ MultiDecoder [5] & 4.1 & 35.7 \\ \hline 50-50 & CycleGAN latent T5 (large) & 16.2 & 60.2 \\ CycleGAN T5 (large) & 69.2 & 94.3 \\ Llama2-7B-Chat & 46.9 & 89.4 \\ \hline CrossAlignment [34] & 3.2 & 62.9 \\ MultiDecoder [5] & 3.2 & 47.4 \\ \hline 66-33 & CycleGAN latent T5 (large) & 15.1 & 46.8 \\ CycleGAN T5 (large) & 69.7 & 96.7 \\ Llama2-7B-Chat & 56.2 & 96.9 \\ \hline CrossAlignment [34] & 2.7 & 60.9 \\ MultiDecoder [5] & 3.2 & 47.4 \\ \hline 66-33 & CycleGAN Istent T5 (large) & 13.9 & 46.9 \\ CrossAlignment [34] & 2.7 & 60.9 \\ MultiDecoder [5] & 2.5 & 49.2 \\ \hline 75-25 & CycleGAN Istent T5 (large) & 13.9 & 46.9 \\ CycleGAN T5 (large) & 69.7 & 96.7 \\ Llama2-7B-Chat & 56.2 & 96.9 \\ \hline 100-0 & CycleGAN Istent T5 (large) & 13.9 & 46.9 \\ CycleGAN T5 (large) & 68.2 & 97.5 \\ Llama2-7B-Chat & 56.1 & 94.1 \\ \hline 00-0 & CycleGAN Istent T5 (large) & 15.5 & 82.9 \\ CycleGAN T5 (large) & 68.2 & 97.5 \\ Llama2-7B-Chat & 56.1 & 94.1 \\ \hline 100-0 & CycleGAN Istent T5 (large) & 15.5 & 82.9 \\ CycleGAN T5 (large) & 92.1 & 94.8 \\ Llama2-7B-Chat & 83.4 & 89.6 \\ \hline \end{array}$	$ \begin{array}{c cccc} CrossAlignment [34] & 8.6 & 49.9 & 20.7 \\ MultiDecoder [5] & 13.2 & 65.4 & 29.4 \\ 0-100 & CycleGAN Iatent T5 (large) & 16.8 & 69.1 & 34.1 \\ CycleGAN T5 (large) & 97.9 & 98.7 & 98.3 \\ Llama2-7B-Chat & 80.4 & 98.9 & 89.2 \\ \hline CrossAlignment [34] & 2.7 & 66.2 & 13.4 \\ MultiDecoder [5] & 2.9 & 54.8 & 12.6 \\ 25-75 & CycleGAN Iatent T5 (large) & 13.9 & 52.9 & 27.1 \\ CycleGAN T5 (large) & 61.1 & 93.9 & 75.7 \\ Llama2-7B-Chat & 52.7 & 97.1 & 71.5 \\ \hline CrossAlignment [34] & 3.2 & 65.3 & 14.5 \\ MultiDecoder [5] & 3.5 & 51.7 & 13.5 \\ 33-66 & CycleGAN Iatent T5 (large) & 16.3 & 51.0 & 28.8 \\ CycleGAN T5 (large) & 65.5 & 94.3 & 78.6 \\ Llama2-7B-Chat & 49.0 & 93.4 & 67.7 \\ \hline CrossAlignment [34] & 3.8 & 63.4 & 15.5 \\ MultiDecoder [5] & 4.1 & 35.7 & 12.1 \\ 50-50 & CycleGAN Iatent T5 (large) & 16.2 & 60.2 & 31.2 \\ CycleGAN T5 (large) & 69.2 & 94.3 & 80.8 \\ Llama2-7B-Chat & 40.9 & 89.4 & 64.8 \\ \hline CrossAlignment [34] & 3.2 & 62.9 & 14.2 \\ MultiDecoder [5] & 3.2 & 47.4 & 12.3 \\ 66-33 & CycleGAN Iatent T5 (large) & 15.1 & 46.8 & 26.6 \\ CycleGAN T5 (large) & 69.7 & 96.7 & 82.1 \\ Llama2-7B-Chat & 56.2 & 96.9 & 73.8 \\ \hline CrossAlignment [34] & 2.7 & 60.9 & 12.8 \\ MultiDecoder [5] & 2.5 & 49.2 & 11.1 \\ 75-25 & CycleGAN Iatent T5 (large) & 13.9 & 46.9 & 25.5 \\ CycleGAN T5 (large) & 69.7 & 96.7 & 82.1 \\ Llama2-7B-Chat & 56.2 & 96.9 & 73.8 \\ \hline CrossAlignment [34] & 2.7 & 60.9 & 12.8 \\ MultiDecoder [5] & 2.5 & 49.2 & 11.1 \\ 75-25 & CycleGAN Iatent T5 (large) & 13.9 & 46.9 & 25.5 \\ CycleGAN T5 (large) & 68.2 & 97.5 & 81.5 \\ Llama2-7B-Chat & 56.1 & 94.1 & 72.7 \\ \hline CrossAlignment [34] & 6.4 & 96.7 & 24.9 \\ MultiDecoder [5] & 9.1 & 78.1 & 26.7 \\ 100-0 & CycleGAN Iatent T5 (large) & 15.5 & 82.9 & 35.8 \\ CycleGAN T5 (large) & 15.5 & 82.9 & 35.8 \\ CycleGAN T5 (large) & 92.1 & 94.8 & 93.4 \\ Llama2-7B-Chat & 83.4 & 89.6 & 86.4 \\ \hline \end{array}$	$ \begin{array}{c} CrossAlignment [34] \\ MultiDecoder [5] \\ 13.2 \\ 65.4 \\ 29.4 \\ 22.0 \\ CycleGAN latent T5 (large) \\ 16.8 \\ 69.1 \\ 34.1 \\ 27.0 \\ CycleGAN T5 (large) \\ 97.9 \\ 98.7 \\ 98.3 \\ 98.4 \\ $	$ \begin{array}{c} \mbox{CrossAlignment [34]} & 8.6 & 49.9 & 20.7 & 14.7 & 52.9 \\ \mbox{MultiDecoder [5]} & 13.2 & 65.4 & 29.4 & 22.0 & 71.3 \\ \mbox{OycleGAN Itatent T5 (large)} & 16.8 & 69.1 & 34.1 & 27.0 & 28.3 \\ \mbox{OycleGAN T5 (large)} & 97.9 & 98.7 & 98.3 & 98.3 & 95.7 \\ \mbox{Llama2-7B-Chat} & 80.4 & 98.9 & 89.2 & 88.7 & 97.1 \\ \mbox{CrossAlignment [34]} & 2.7 & 66.2 & 13.4 & 5.2 & 61.2 \\ \mbox{MultiDecoder [5]} & 2.9 & 54.8 & 12.6 & 5.5 & 55.3 \\ \mbox{CycleGAN 15 (large)} & 61.1 & 93.9 & 75.7 & 74.0 & 95.6 \\ \mbox{Llama2-7B-Chat} & 52.7 & 97.1 & 71.5 & 68.3 & 92.0 \\ \mbox{CrossAlignment [34]} & 3.2 & 65.3 & 14.5 & 6.1 & 61.1 \\ \mbox{MultiDecoder [5]} & 3.5 & 51.7 & 13.5 & 6.6 & 51.9 \\ \mbox{CycleGAN 15 (large)} & 65.5 & 94.3 & 78.6 & 77.3 & 95.4 \\ \mbox{Llama2-7B-Chat} & 49.0 & 93.4 & 67.7 & 64.3 & 89.7 \\ \mbox{CycleGAN 15 (large)} & 69.2 & 94.3 & 80.8 & 79.8 & 93.0 \\ \mbox{Llama2-7B-Chat} & 49.0 & 93.4 & 67.7 & 64.3 & 89.7 \\ \mbox{CrossAlignment [34]} & 3.2 & 62.9 & 14.2 & 6.1 & 60.5 \\ \mbox{MultiDecoder [5]} & 16.2 & 60.2 & 31.2 & 25.5 & 53.6 \\ \mbox{CycleGAN 15 (large)} & 69.2 & 94.3 & 80.8 & 79.8 & 93.0 \\ \mbox{Llama2-7B-Chat} & 46.9 & 89.4 & 64.8 & 61.5 & 86.2 \\ \mbox{CycleGAN 15 (large)} & 69.7 & 96.7 & 82.1 & 81.0 & 95.2 \\ \mbox{Llama2-7B-Chat} & 56.2 & 96.9 & 73.8 & 71.1 & 94.0 \\ \mbox{CrossAlignment [34]} & 3.2 & 62.9 & 14.2 & 6.1 & 60.5 \\ \mbox{MultiDecoder [5]} & 2.5 & 49.2 & 11.1 & 4.8 & 52.4 \\ \mbox{CycleGAN 15 (large)} & 69.7 & 96.7 & 82.1 & 81.0 & 95.2 \\ \mbox{Llama2-7B-Chat} & 56.2 & 96.9 & 73.8 & 71.1 & 94.0 \\ \mbox{CrossAlignment [34]} & 2.7 & 60.9 & 12.8 & 5.2 & 57.8 \\ \mbox{MultiDecoder [5]} & 2.5 & 49.2 & 11.1 & 4.8 & 52.4 \\ \mbox{CycleGAN 15 (large)} & 69.7 & 96.7 & 82.1 & 81.0 & 95.2 \\ \mbox{Llama2-7B-Chat} & 56.2 & 96.9 & 73.8 & 71.1 & 94.0 \\ \mbox{CrossAlignment [34]} & 2.7 & 60.9 & 12.8 & 5.2 & 57.8 \\ \mbox{MultiDecoder [5]} & 2.5 & 49.2 & 11.1 & 4.8 & 52.4 \\ \mbox{CycleGAN 15 (large)} & 68.2 & 97.5 & 81.5 & 80.3 & 94.6 \\ \mbox{Llama2-7B-Chat} & 56.1 & 94.1 & 72.7 & 70.3 & 91.5 \\ C$	CrossAlignment [34] 8.6 49.9 20.7 14.7 52.9 21.3 MultiDecoder [5] 13.2 65.4 29.4 22.0 71.3 30.7 0-100 CycleGAN Itsent T5 (large) 16.8 69.1 34.1 27.0 28.3 21.8 CycleGAN T5 (large) 97.9 98.7 98.3 98.3 95.7 96.8 Llama2-7B-Chat 80.4 98.9 89.2 88.7 97.1 88.4 CrossAlignment [34] 2.7 66.2 13.4 5.2 61.2 12.9 MultiDecoder [5] 2.9 54.8 12.6 55.5 55.3 12.7 25-75 CycleGAN Itsent T5 (large) 61.1 93.9 75.7 74.0 95.6 76.4 Llama2-7B-Chat 52.7 97.1 71.5 68.3 92.0 69.6 CycleGAN Itsent T5 (large) 16.3 51.0 28.8 24.7 49.8 28.5 CycleGAN Itsent T5 (large) 16.2 63.4 15.5

In our qualitative analysis, we identify the main failure cases and summarize them below:

- *Metaphoric language:* the model's limited ability to accurately recognize and modify metaphoric language and idiomatic expressions. This can result in the model retaining the original expressions in the generated text, which may not conform to the desired style.
- *Slang:* the model's difficulty in accurately recognizing and modifying common words used with their slang meaning, particularly in cases where the conversion is from informal to formal language. In such cases, the model may consider the slang word as already formal and fail to convert it, resulting in the retention of the original expression.

Table 16. Effect of the λ_{cyc} hyperparameter without classifier-guided loss on the GYAFC-family dataset with BART (base) model - ref-BLEU (ref-B avg), style accuracy with BERT and TextCNN (acc_{BERT}, acc_{CNN}), geometric mean (GM) and harmonic mean (HM) of ref-BLEU and style accuracy.

λ_{cyc}	ref-B avg	acc_{BERT}	GM	HM	acc _{CNN}	GM	HM
0	39.0	6.1	15.4	10.6	6.8	16.3	11.6
0.1	43.2	46.2	44.7	44.7	45.2	44.2	44.2
1	43.1	47.6	45.3	45.2	46.5	44.8	44.7
10	42.8	44.0	43.4	43.4	43.1	42.9	42.9
50	43.5	47.4	45.4	45.4	46.6	45.0	45.0
100	43.5	50.1	46.7	46.6	49.6	46.5	46.4

Table 17. Effect of the classifier-guided loss on the GYAFC-family dataset - ref-BLEU (ref-B avg), style accuracy with BERT and TextCNN (accBERT, accCNN), geometric mean (GM) and harmonic mean (HM) of ref-BLEU and style accuracy.

1 407								
1437		ref-B avg	acc _{BERT}	GM	HM	acc _{CNN}	GM	HM
1438	CycleGAN BART (base)	43.7	50.7	47.1	46.9	49.4	46.5	46.4
1439	w/o style classifier	42.8	44.0	43.4	43.4	43.1	42.9	42.9
1440	CycleGAN BART (large)	43.5	50.8	47.0	46.9	49.9	46.6	46.5
1441	w/o style classifier	43.4	47.9	45.6	45.5	47.6	45.5	45.4
1442	CycleGAN T5 (small)	42.1	38.7	40.4	40.3	39.6	40.8	40.8
1443	w/o style classifier	42.1	39.2	40.6	40.6	39.9	41.0	41.0
1444	CycleGAN T5 (base)	44.0	47.7	45.8	45.8	46.2	45.1	45.1
1445	w/o style classifier	44.0	47.1	45.5	45.5	45.7	44.8	44.8
1440	CycleGAN T5 (large)	45.4	59.5	52.0	51.5	58.1	51.4	51.0
1447	w/o style classifier	44.4	49.2	46.7	46.7	47.9	46.1	46.1
1110	-	1						

Table 18. Effect of the classifier-guided loss on the GYAFC-family dataset | informal \rightarrow formal – ref-BLEU (ref-B), style accuracy with BERT and TextCNN (accBERT, accCNN), geometric mean (GM) and harmonic mean (HM) of ref-BLEU and style accuracy.

1453		ref-B	acc _{BERT}	GM	HM	acc _{CNN}	GM	HM
1454	CycleGAN BART (base)	58.7	42.2	49.8	49.1	47.0	52.5	52.2
1455	w/o style classifier	57.2	34.5	44.4	43.0	38.4	46.9	46.0
1456	CycleGAN BART (large)	59.3	41.2	49.4	48.6	46.2	52.3	51.9
1457	w/o style classifier	59.9	47.8	53.5	53.2	54.4	57.1	57.0
1458	CycleGAN T5 (small)	54.3	28.7	39.5	37.6	34.8	43.5	42.4
1459	w/o style classifier	54.3	29.1	39.8	37.9	35.3	43.8	42.8
1460	CycleGAN T5 (base)	56.7	28.9	40.5	38.3	32.1	42.7	41.0
1461	w/o style classifier	57.7	32.3	43.2	41.4	36.3	45.8	44.6
1402	CycleGAN T5 (large)	61.9	49.2	55.2	54.8	53.2	57.4	57.2
1465	w/o style classifier	58.0	32.3	43.3	41.5	37.2	46.4	45.3

• In-depth rephrasing: the model's inability to perform more profound rephrasing of the input sentence when necessary to achieve the desired style transfer. This is a common limitation of non-parallel TST approaches, where the lack of parallel training data makes it challenging to learn more complex mappings between styles.

ACM Trans. Intell. Syst. Technol., Vol. 1, No. 1, Article . Publication date: July 2024.

2	1
- 5	
0	

1471	Table 19. Effect of the classifier-guided loss on the GYAFC-music dataset - ref-BLEU (ref-B avg), style
1472	accuracy with BERT and TextCNN (accBERT, accCNN), geometric mean (GM) and harmonic mean (HM) of
1473	ref-BLEU and style accuracy. * denotes results from the paper.

1474								
1475		ref-B avg	acc _{BERT}	GM	HM	acc _{CNN}	GM	HM
1476	CycleGAN BART (base)	43.6	57.2	49.9	49.5	57.8	50.2	49.7
1477	w/o style classifier	42.3	49.0	45.5	45.4	50.6	46.3	46.1
1478	CycleGAN BART (large)	42.0	43.1	42.5	42.5	43.8	42.9	42.9
1479	w/o style classifier	40.8	41.8	41.3	41.3	43.0	41.9	41.9
1480	CycleGAN T5 (small)	40.6	37.9	39.2	39.2	39.0	39.8	39.8
1481	w/o style classifier	40.6	37.6	39.1	39.0	38.6	39.6	39.6
1482	CycleGAN T5 (base)	42.0	45.4	43.7	43.6	47.7	44.8	44.7
1483	w/o style classifier	42.0	43.2	42.6	42.6	45.5	43.7	43.7
1484	CycleGAN T5 (large)	45.6	70.5	56.7	55.4	70.1	56.5	55.3
1485	w/o style classifier	43.3	43.0	43.1	43.1	45.1	44.2	44.2

1486
1487Table 20. Effect of the classifier-guided loss on the GYAFC-music dataset | informal \rightarrow formal - ref-BLEU
(ref-B avg), style accuracy with BERT and TextCNN (acc_{BERT}, acc_{CNN}), geometric mean (GM) and harmonic
mean (HM) of ref-BLEU and style accuracy.1489

1490		ref-B	acc _{BERT}	GM	HM	acc _{CNN}	GM	HM
1491	CycleGAN BART (base)	57.6	53.7	55.6	55.6	48.4	52.8	52.6
1492	w/o style classifier	55.5	43.8	49.3	49.0	39.8	47.0	46.4
1493	CycleGAN BART (large)	55.3	42.8	48.7	48.3	38.4	46.1	45.3
1494	w/o style classifier	52.3	37.2	44.1	43.5	32.9	41.5	40.4
1495	CycleGAN T5 (small)	51.6	34.4	42.1	41.3	29.5	39.0	37.5
1496	w/o style classifier	51.7	35.5	42.8	42.1	29.8	39.3	37.8
1497	CycleGAN T5 (base)	55.2	36.5	44.9	43.9	33.4	42.9	41.6
1498	w/o style classifier	54.9	35.7	44.3	43.3	32.5	42.2	40.8
1499	CycleGAN T5 (large)	62.7	67.1	64.9	64.8	61.6	62.1	62.1
1500	w/o style classifier	58.3	48.5	53.2	53.0	44.4	50.9	50.4
1501		1				1		

Table 21. Effect of the classifier-guided loss on the Yelp dataset – ref-BLEU (ref-B avg), style accuracy with
 BERT and TextCNN (acc_{BERT}, acc_{CNN}), geometric mean (GM) and harmonic mean (HM) of ref-BLEU and
 style accuracy.

1506		ref-B avg	acc_{BERT}	GM	HM	acc _{CNN}	GM	HM
1507	CycleGAN BART (base)	55.7	78.8	66.3	65.3	77.8	65.8	64.9
1508	w/o style classifier	53.1	82.7	66.3	64.7	81.3	65.7	64.2
1509	CycleGAN BART (large)	56.5	75.1	65.1	64.5	74.6	64.9	64.3
1510	w/o style classifier	56.9	73.9	64.8	64.3	73.1	64.5	64.0
1511	CycleGAN T5 (small)	53.0	78.0	64.3	63.1	78.2	64.4	63.2
1512	w/o style classifier	54.4	76.5	64.5	63.6	77.8	65.1	64.0
1513	CycleGAN T5 (base)	54.2	76.6	64.4	63.5	77.3	64.7	63.7
1514	w/o style classifier	55.4	74.2	64.1	63.4	74.7	64.3	63.6
1515	CycleGAN T5 (large)	55.3	72.9	63.5	62.9	73.7	63.8	63.2
1516	w/o style classifier	56.6	71.9	63.8	63.3	71.9	63.8	63.3
1517								

1517

1505

1520Table 22. Effect of different pre-trained style classifier models on the GYAFC-family dataset with BART (base)1521model - ref-BLEU (ref-B avg), style accuracy with BERT and TextCNN (acc_{BERT}, acc_{CNN}), geometric mean1522(GM) and harmonic mean (HM) of ref-BLEU and style accuracy.

style classifier	ref-B avg	acc _{BERT}	GM	HM	acc _{CNN}	GM	HM
BERT-base	43.7	50.7	47.1	46.9	49.4	46.5	46.4
BERT-large	43.6	48.2	45.8	45.8	47.7	45.6	45.6
RoBERTa-base	43.6	49.1	46.3	46.2	47.9	45.7	45.7
RoBERTa-large	43.6	50.3	46.8	46.7	49.2	46.3	46.2
DistilBERT-base	42.7	46.1	44.4	44.3	44.7	43.7	43.7

1529 1530

1531 1532

1533

A more detailed analysis of the qualitative results follows.

1534 F.1 Formality transfer

Comparison with the reference annotations. Tables 23 and 24 include success and failure cases for 1535 both directions of the formality transfer task, for both the family and music domains, respectively. 1536 Considering the family domain, the first two examples in Table 23 shows the ability of the model 1537 to transform an informal text into the corresponding formal version. Specifically, in the first case 1538 the model capitalizes the first letter of the sentence, the contracted and colloquial form "dont" 1539 is converted into "do not", the full stop is added at the end of the text and the slang terms "dat" 1540 and "da" are transformed into their proper versions "that" and "the", respectively. Although two 1541 occurrences of the word "dat" are present in the input document, to avoid repetitions the model 1542 does not replace both of them with "that" but one of the two is mapped to "who", thus denoting 1543 both language variety in the generated text and the model's understanding capabilities to recognize 1544 nuanced differences in language use. In the second case, the model correctly introduces the proper 1545 subject in the output sentence and replaces the smiling emoticon with a full stop, resulting in a 1546 perfect match with one of the available references. Considering the formal-to-informal direction, 1547 the model transforms the first letter of the sentence into its lowercase version, replaces the subject 1548 you" with its contracted form "u" and rephrases the wording "in that manner" with the more 1549 informal version "like that". In the second example, similar transformations are applied by the 1550 model (i.e., lowercase for the first letter, "you" \rightarrow "u"); moreover, "do not" is contracted into "dont", 1551 "are" is replaced by the informal abbreviation "r" and punctuation is removed from the sentence. 1552 Although the generated text does not match any of the references, it is successfully transferred to 1553 the informal style. 1554

Considering now the potential failure cases reported in the bottom part of Table 23, we can 1555 observe that in one example the model successfully capitalizes the first letter and adds a full stop at 1556 the end of the sentence, but fails to modify the informal language in the remaining text. This is 1557 probably caused by the use of the metaphoric expression "hit the nail on the head" which the model 1558 may not recognize as informal language and, therefore, fails to modify. In the second example, 1559 the model successfully applies some modifications to improve the formality of the input text, 1560 such as capitalizing the first letter and replacing "ur" with "you are". However, the model fails 1561 to recognize the slang term "hoe" and, therefore, copies it to the output text, likely because it is 1562 interpreted in its literal meaning. This example highlights a limitation of the model in accurately 1563 recognizing and modifying slang words, which can result in the retention of inappropriate language 1564 in the generated text. Considering both failure cases in the formal-to-informal direction, the model 1565 introduces informal elements in the text (i.e., lowercase, contractions, punctuation removal) but 1566 the generated sentences do not match the corresponding references. This limitation in accurately 1567

¹⁵⁶⁸

modifying informal language may be due to the examples' complexity, where significant rephrasing
is required to achieve the desired style transfer. Such cases can be particularly challenging for
non-parallel TST approaches, where the lack of parallel training data makes it difficult to learn
complex mappings between styles.

Table 24 reports a selection of success and failure cases for both directions of the formality transfer 1573 task within the music domain. In the first two examples, the model correctly transforms the informal 1574 samples into their formal counterparts. More precisely, in the first case several modifications are 1575 applied: the first letter is capitalized, abbreviations are mapped to their extended versions (i.e., 1576 $2^{"} \rightarrow$ "to", "ur" \rightarrow "your", "u" \rightarrow "you") and three exclamation marks are replaced with a full 1577 stop. In the second example, in addition to similar modifications performed in previous cases, the 1578 model corrects the spelling of the word "like"; moreover, the generated text adopts the proper 1579 capitalization while the source sentence is entirely written in uppercase. This results in an almost 1580 perfect match with one of the proposed references. Considering the formal-to-informal direction, 1581 the model correctly uses contractions, abbreviations and no capitalization when rewriting the 1582 input text in informal style. From the second example, it is also possible to see that the model 1583 replaces the words "recall" and "television" with the more common alternatives "remember" and 1584 "tv", respectively. 1585

In the first failure case, the model tries to convert the source text into its corresponding formal 1586 version by capitalizing the first letter and adding a punctuation mark at the end of the sentence but 1587 it is not able to correct and substitute the words "no" and "sight" with their homophones "know" 1588 and "site", respectively. In the second example, the generated text closely resembles a copy of 1589 the input sentence. This can be attributed to the unconventional formatting of the word "respect" 1590 which is written in uppercase letters with each letter separated by a dot. Such formatting may 1591 have caused the model to interpret it as an acronym or a specific entity, making it challenging 1592 to effectively transform and generate the desired output. In the first failure case of the informal-1593 to-formal direction, the generated sentence displays certain features commonly seen in informal 1594 texts, such as the absence of a subject and contracted forms. However, it does not correspond to 1595 any of the provided references. It is worth noting that the model fails to recognize and modify the 1596 expression "under the weather", which likely contributes to the discrepancy between the generated 1597 output and the desired reference sentences. It is possible that the model did not understand the 1598 idiom "under the weather", which typically refers to someone feeling unwell or sick. Finally, in the 1599 last example, the model nearly replicates the input text without substantial modifications. Once 1600 again, a more elaborate rewriting would be necessary to achieve a closer match with the proposed 1601 references. 1602

Comparison with the baseline methods. Tables 25 and 26 report several examples of generated 1604 outputs for the formality transfer task for both family and music domains, respectively. All baseline 1605 methods struggle with preserving the source text content while converting style attributes from 1606 informal to formal and vice versa (see Table 25). Conversely, our approach achieves better preser-1607 vation of the original content. According to the quantitative results reported in Section 5.5, the 1608 formality transfer task is more challenging than the sentiment transfer task, especially for content 1609 preservation. As an example, in the second case of the informal-to-formal direction from the family 1610 domain, only two out of five baselines include the word "avatar" in the generated sentence while 1611 the outputs of the remaining three practically have no overlap with the source text. On the contrary, 1612 our model perfectly matches one of the proposed references. By incorporating cycle-consistency 1613 losses, the model is encouraged to maintain a certain level of coherence and consistency in the 1614 generated text while effectively transforming style attributes. This aspect contributes to the better 1615 preservation of the original content, addressing a common challenge faced by baseline methods. 1616 1617

1619						
1620					$informal \rightarrow formal$	$formal \rightarrow informal$
1020			Sourc	ce	i dont like a guy dat does dat all da time	Do you feel in that manner?
1621			Gene	rated	I do not like a guy who does that all the time.	do u feel like that?
1622					I do not like men that does that all of the time.	do u fell that way ?
1(00			Refer	ences	I do not find a man appealing who behaves that way.	now u Ieel Is that how you fael?
1025					I do not like a guy that does that all of the time	is that how u feel?
1624	SI	iccess	Sourc	re	Sounds like a rhetorical question :)	I do not know perhaps you are young
1625			Gene	rated	It sounds like a rhetorical question.	i dont know perhaps u r young
					It sounds like a rhetorical question.	Maybe your young, I don't know.
1626			Dafan		That sounds like a rhetorical question.	maybe your young IDK.
1627			Refer	ences	That sounds like a rhetorical question.	MAYBE AGE IS AN ISSUE
1628					Are you posing a rhetorical question?	maybe you're just young, idk
1020			Sourc	ce	bembie hit the nail on the head	I can lift him.
1629			Gene	rated	Bembie hit the nail on the head.	I can lift him :)
1630					You are speaking the truth.	yup, i can lift him, he's only 185lol
1 (0 1			Refer	ences	Bembie got the answer correct.	i can lift him no problem
1631			licici	ciices	Bembie got it right.	Yep I can lift him.
1632	f	ailure			Bembie reached the proper conclusion.	He's a weakling, I can pick him up!
1633			Sourc	ce	so in my opinon ur not a hoe!	It is a very long list.
1055			Gene	rated	So in my opinion you are not a noe:	Its a very long list
1634					In my opinion you are not a hoe	its a long ass list
1635			Refer	ences	So in my opinion you are not a promiscuous woman	That list could go on forever
1626					In my opinion, you are not promisedous woman.	the list is endless, thh
1050					, , , , , , , , , , , , , , , , , , ,	
1638 1639			Tab		informal → formal	formal \rightarrow informal
1640		Sourc	e	explai	n 2 ur parents that u really want 2 act !!!	Can you not rhyme the words together?
1641		Gener	rated	Expla	in to your parents that you really want to act.	cant u rhyme the words together?
1642				Expla	in to your parents that you want to act.	Cant you just rhyme the words with themselves?
1012		Refer	ences	Expla	in to your parents that you very much want to act.	Can ya not rhyme da damn wo ds togeder?
1643				Evola	in to your parents that you truly want to act	can't you just rhyme the words???
1644	success	Sourc	e	LLKE	THEM ALSO THEY ARE SOOOO CUTE!	I cannot recall the last time I turned on my television
1645		Gener	rated	I like	them. Also, they are so cute!	i cant remember the last time i turned on my tv
1045				I like	them; also, they are adorable.	I can't remember when I last turned on my TV.
1646		Defer	-	I like	them, also. They are so cute.	I don't remember when I turned my TV on.
1647		Refer	ences	I like	them because they are cute.	i dont remember when i watched tv last
1749				I like	them as well, they are so cute!	I dont even remember the last time i watched tv
1048		Sourc	e	does a	nybody no a web sight for a casting call	I am unsure, as I am currently under the weather.
1649		Gener	rated	Does	anybody no a web sight for a casting call?	unsure, as im currently under the weather.
1650				Does	anyone know a website for a casting call?	Umm, don't know, my head hurts
1/51		Refer	ences	what	website has casting calls?	I dunno, I m sick, dude
1051				Is any	one aware of a Casting Call website?	i don t know as i m feeling sick.
1652	failure	Sourc	e	Heim	st needs to RESPECT me my pplz_lol	L can't listen
1653		Gener	rated	He ju	st needs to R.E.S.P.E.C.T. me my pplz. Iol	I cant listen.
1/54				He ne	eds to respect me and my people.	I turned my ears away !
1654		D		He ne	eds to show respect to me and my friends.	Not listening
1655		Keter	ences	He sh	ould respect myself and my colleagues.	I didn't want to hear it
1656				he jus	t need to respect me and my people.	La la la I can't hear you!
1/57						
100/						

Table 23. Examples from GYAFC-family test set - success and failure cases.

DeepLatent [7] is the only baseline in the music domain whose outputs are publicly available. In almost all four reported examples, it generates sentences significantly different from the source text (see Table 26). This results in poor performance in the preservation of the original meaning, as confirmed also by the much lower BLEU scores with respect to our model.

F.2 Sentiment transfer

Comparison with the reference annotations. We analyze sentiment transfer cases in both directions (see Table 27). Considering the first negative-to-positive success case, the model not only removes

1665 1666

1658

1659

1660

1661 1662

1663

1664

ACM Trans. Intell. Syst. Technol., Vol. 1, No. 1, Article . Publication date: July 2024.

Table 25. Examples from GYAFC-family test set - comparison with baselines.

	$informal \rightarrow formal$	$formal \rightarrow informal$
Source	i guess its normal, but id just move on	I do not think you are ready if you have to ask other people.
	I suppose it is normal, but I would just move on.	if u have to ask other people i really dont think ur ready to
0 - f	I imagine it's common, however, I'd simply proceed.	If u ask others, I don't think ur ready to
Aererences	I guess that it is normal, but I would move on.	you ain ready if u gotta ask the folks
	I guess it is normal but I would just move on.	youre not ready if u need to ask others
(etrieveOnly [18]	it typically on from there .	but i 'm and think thats - maybe i am a .
CrossAlignment [34]	i would not care about them and they are attractive .	i would tell him if you just want to get out with him
InsupervisedMT [45]	i am not normal , but id just move on to the club .	i do n't think you are ready if you have to ask other people it is .
DualRL [23]	i guess it is normal , but just move on .	i dont think you are ready if you have to ask other people .
NASTLatentLearn [11]	i guess it normal , but i just do on	i do not think you are ready if you have to ask other people .
Ours	I guess it is normal, but I would just move on.	i dont think ur ready if u have to ask other people
ource	r u talking about ur avatar?	I apologize, but I do not know. I wish the both of you luck.
	Are you referring to your avatar?	hey im sorry i don't know and i wish best of luck to both of you
aferences	Are you talking about your avatar?	Im sorry, idk.i wish u both luck
elefences	Are you talking about your avatar?	sorry i dont know, good luck to you both!
	Are you talking about your avatar?	sorry, but idk. I still wish you the best.
RetrieveOnly [18]	it is not pleasant driving a - 100 .	yes , your on the wrong .
CrossAlignment [34]	you should be happy with marriage .	i think about the same thing i would know what you like me .
JnsupervisedMT [45]	near you talking about passion , without even if it .	i apologize , but i do n't know i wish you the both of you hookin
DualRL [23]	it is talking about avatar avatar ?	i dunno er
NASTLatentLearn [11]	do you talking about it avatar ,	i do , but i do not know i wish the both of you luck .
Durs	Are you talking about your avatar?	sorry but i dont know i wish the both of you luck
Table	26. Examples from GYAFC-music test set	– comparison with baselines.
	•	-

(00		$informal \rightarrow formal$	$formal \rightarrow informal$
688	Source	im a huge green day fan!!!!!	How old are you? You should be at least 18 years old.
689		I am a huge fan of the band Green Day.	i dont know how old ru?ur supposed to be 18
(00	Defense	I am a big fan of the band Green Day.	ur supposed to be 18, how old r u?
690	References	I am a big fan of Green Day!	Are you over 18? How old r u?
691		I am a huge Green Day fan.	how old is ur age? u must be at least 18
602	DeepLatent [7]	I am a fan of the movie .	How old are you you should be at least 18 years old .
092	Ours	I am a huge Green Day fan.	how old r u u should be at least 18 years old
693	Source	YOUR WASTING YOUR TIME, SAYS THE BOY.	It is not close to a PlayStation Portable (PSP), but it is close enough.
694		The boy said you are wasting your time.	Not close to a PSP, but close enough.
COF	Deferences	The boy says you're wasting your time.	ITS NOT A PSP BUT ITS YOURS
095	References	The boys said, "You are wasting your time".	Eh its not exactly a PSP but its good enough
696		"You are wasting your time", says the boy.	it isn't close to a PSP but its close enuf
697	DeepLatent [7]	Try the fourth of Narnia , but I am not sure .	i think its going to get a goddamn door (but it is easy to get enough .
	Ours	"You are wasting your time", says the boy.	its not close to a psp but its close enough

1698 1699

1

1667

the negation "not" but also strengthens the positive sentiment by adding the word "definitely", 1700 as done by the first reference. In the second case, the model is able to modify two aspects of the 1701 input sentence, therefore confirming its ability to deal with different aspects at the same time. 1702 Considering the positive-to-negative direction, the first example shows that the model correctly 1703 modifies both negative adjectives converting them into their positive counterparts. According to 1704 the second example, it can be seen that, in addition to modifying the sentiment of the input text, 1705 the model shows solid language understanding capabilities: instead of simply adding "limited", it 1706 also replaces "variety" with "number". Although this may be suboptimal in terms of BLEU score, 1707 language diversification is generally a positive feature. In all the reported success cases, it can be 1708 noticed that the model not only correctly modifies the sentiment of the source sentence but does it 1709 in a consistent way from a semantic point of view (e.g., over fried \rightarrow perfectly fried, hard \rightarrow soft, 1710 amazing \rightarrow disgusting). 1711

Moving to the failure cases, in the first example, the model almost copies the input sentence by adding the word "good" but the expression "run of the mill" is kept unchanged. The same pattern can be observed in the positive-to-negative case where the source sentence is copied to the output

1717				
1710			negative \rightarrow positive	$positive \rightarrow negative$
1/18		Source	we will not be coming back .	food is very fresh and amazing !
1719		Generated	we will definitely be coming back.	food is very bland and disgusting!
1720			we will definitely be coming back !	food was old and stale
1720		Peferences	we will be back soon,	food is not fresh and average
1721		References	we will come back soon	the food is very bad and disgusting
1722	511000655		we will quickly come back	the food is not good and disgusting !
1722	success	Source	it was over fried and very hard .	the variety of sushi rolls makes for a good eating .
1723		Generated	it was perfectly fried and very soft.	the limited number of sushi rolls makes for a bad eating.
1794			it was a little overcooked	there is limited variety for sushi rolls .
1724		Peferences	it was cooked at its point and very soft , exquisite	the variety of sushi roll does not tasted good
1725		References	it was not over fried and not hard,	no variety of sushi rolls makes for a no good eating .
1726			it was fried just right and very tender	the variety of sushi rolls make for a bad eating
1720		Source	even the ice cream was very , very run of the mill .	the soup is absolutely to die for .
1727		Generated	even the ice cream was very good, very run of the mill.	the soup is absolutely to die for.
1728			even the ice cream was very , very unique and good .	the soup will make you want to die
		Deferences	the ice cream is very nice .	the soup is dreadful .
1729		References	the ice cream is good	the soup is bad .
1730	failura		the ice cream is very perfect .	the soup is not good
	lanure	Source	what about the time i wasted ?	bartender made a great dirty martini !
1731		Generated	what about the time i saved?	bartender made a dirty martini!
1732			service was quick and swift	bartender gave me a watered down martini
		Peferences	i really enjoy my time	bartender made a bad martini
1733		increnences	this is really worth it	bartener made a nasty dirty martini
1734			i enjoy my time .	bartender made a bad dirty martini !

Table 27. Examples from Yelp test set – success and failure cases.

1736 without modification. This is probably due to the use of the expression "to die for", which may be 1737 erroneously recognized as already negative and therefore is not modified. Similarly to the formality 1738 transfer task, these examples highlight a challenge faced by our model when dealing with texts 1739 that include idiomatic expressions. The model struggles to correctly identify and handle these 1740 expressions, leading to failure cases where the sentiment transfer is not accurately achieved. In the 1741 second negative-to-positive failure case, the model attempted to transform the word "wasted" into 1742 "saved", resulting in a shift towards a positive nuance. However, it is important to note that this 1743 generated output does not align with any of the provided references. The discrepancy between 1744 the model's output and the references can be attributed to the fact that the desired sentiment 1745 expression in the references may require a more significant rephrasing of the input sentence. 1746 Since our approach is trained in a self-supervised setting, where explicit supervision for specific 1747 rephrasing patterns is not provided, achieving a closer match to the references in such cases 1748 becomes more challenging. Lastly, the final example demonstrates another limitation of the model, 1749 where it fails to recognize the term "dirty martini" and, similar to a previous case, incorrectly 1750 assumes a negative sentiment in the sentence. Consequently, the model does not explicitly modify 1751 the sentiment of the text but only removes the adjective "great". This indicates that the model's 1752 performance is hindered when encountering domain-specific terms or expressions that are not 1753 adequately identified and processed.

Comparison with the baseline methods. When comparing our approach with the baseline methods, 1755 it is observed that the majority of the baselines achieve successful sentiment transfer. However, 1756 similarly to the formality transfer task, they struggle to preserve the original content of the text, as 1757 indicated in Table 28. Let us consider, as an example, the second case reported for the negative-to-1758 positive direction. Our model perfectly matches one of the available references whereas the baselines 1759 either do not transfer the style (e.g., DualRL [23]) or modify the text content (e.g., GTAE [35], 1760 NASTLatentLearn [11]). Considering the first example in the positive-to-negative case, all the 1761 baselines considered keep the adjective "hot" unchanged and two of them also do not modify 1762 the positive adverb "perfectly", whereas our model correctly transfers the sentiment of the entire 1763 1764

36

1716

1735

1754

ACM Trans. Intell. Syst. Technol., Vol. 1, No. 1, Article . Publication date: July 2024.

1765

Table 28. Examples from Yelp test set - comparison with baselines.

1766			
1767	-0	$negative \rightarrow positive$	$positive \rightarrow negative$
1769	Source	if i could give zero stars i def would.	it's hot, cooked perfectly, and delicious !
1708		if i could give 5+ stors and is great	it's cold cooked imperfectly and bad taste
1769	References	if i could give ten stars, i would definitely do it	it's cold, cooked unperfectly, and suck !
1770		if i could give 5+ stars i def would	it 's cold , not cooked perfectly , and taste bad
1771	RetrieveOnly [18]	best part is , everything is made from scratch and you .	maybe the hot dog is cold , but the chili is hot .
1//1	DualRL [23]	if i could give it a def would def recommend it .	it 's hot , over cooked , and cold !
1772	GTAE [35]	if i could give perfect stars i def deliciously .	it 's hot , cooked terrible , and disappointing !
1773	NASILatentLearn [11]	if i could give a stars i def would .	it's hot, cooked perfectly, and bland !
1774	Ours	if i could give ten stars i def would	it's cold overcooked and bland!
	Causa	tested nelles ald i send n't balines it	this place is supercontent !
1775	Source	tasted really fresh i could n't believe it	this place is super vucky !
1776		tasted really fresh, i could n't believe it	this place is not vummy at all !
1777	References	tasted really new , i could n't believe it	this place is terrible !
1778		very new taste , it is believable	this place is lacking in taste
1770	RetrieveOnly [18]	really really really strong margaritas !	i would give this restaurant a zero , if that was an option .
1779	DualRL [23]	tasted really old, i could definitely believe it.	this place is super yummy ?
1780	GIAE [35] NASTI atenti earn [11]	i really good i could it believe it	this place is terribly boring !
1781	MixAndMatch [24]	tasted really amazing, i could n't believe it.	this place is so murky !
1790	Ours	tasted really fresh, i couldn't believe it.	this place is super yuck!
1782			
1783			
1784	sentence. Finally, i	in the last example, our model is the only	one that properly converts "yummy"
1785	into "vuck". This co	onfirms the ability of our model to transfer	the style consistently from a semantic
1797	point of view as a	beerved in the analysis of success cases	, ,
1700	point of view, as o	bserved in the analysis of success cases.	
1787			
1788			
1789			
1700			
1790			
1791			
1792			
1793			
1704			
1794			
1795			
1796			
1797			
4700			
1798			
1799			
1800			
1801			
1000			
1802			
1803			
1804			
1805			
1907			
1900			
1807			
1808			
1809			
1010			
1810			
1811			
1812			
1813			
1015			