

Style Transfer with Multi-iteration Preference Optimization

Anonymous ACL submission

Abstract

Numerous recent techniques for text style transfer characterize their approaches as variants of reinforcement learning and preference optimization. In this work, we consider the relationship between these approaches and a class of optimization approaches developed primarily for (non-neural) statistical machine translation, formerly known as ‘tuning’. Inspired by these techniques from the past, we improve upon established preference optimization approaches, incorporating multiple iterations of exploration and optimization, and choosing contrastive examples by following a ‘hope’ vs ‘fear’ sampling strategy. Cognizant of the difference between machine translation and style transfer, however, we further tailor our framework with a new pseudo-parallel generation method and a dynamic weighted reward aggregation method to tackle the lack of parallel data and the need for a multi-objective reward. We evaluate our model on two commonly used text style transfer datasets. Through automatic and human evaluation results we show the effectiveness and the superiority of our model compared to state-of-the-art baselines.

1 Introduction

Text style transfer aims to rewrite a given text to match a specific target style while preserving the original meaning. This task has drawn significant attention recently due to its broad range of applications, such as text simplification (Laban et al., 2021), formality transfer (Rao and Tetreault, 2018; Liu et al., 2022), text detoxification (Dale et al., 2021; Hallinan et al., 2023b), authorship transfer (Patel et al., 2023; Liu et al., 2024), and authorship anonymization (Shetty et al., 2018; Bo et al., 2021). Recent approaches have focused on pseudo-parallel data generation (Krishna et al., 2020; Riley et al., 2021) and policy optimization (Gong et al., 2019; Liu et al., 2021b). STEER (Hallinan et al.,

2023a) and ASTRAPOP (Liu et al., 2024) combine the two and achieve state-of-the-art performance on text style transfer and authorship style transfer, respectively.

In this work, we seek to advance the frontier of text style transfer, drawing inspiration from the optimization techniques developed in the era of statistical phrasal machine translation, in which the lack of correlation between the log-linear model objective and the desired evaluation metric, typically BLEU (Papineni et al., 2002), was observed (Och, 2003). Approaches to align¹ the two objectives came to be known as *tuning*,² beginning with Och (2003), and evolving into online variants (Chiang et al., 2008), rank-based approaches (Hopkins and May, 2011), batch-based approaches (Cherry and Foster, 2012), and several others. Tuning methods follow a generate-and-optimize pattern: a model is used to generate multiple candidate hypotheses per input, and then parameters are adjusted such that the argmax according to the model score also maximizes the evaluation metric. In this regard, tuning methods resemble approaches taken in the application of policy optimization algorithms, such as PPO (Schulman et al., 2017), to generative language modeling (Ouyang et al., 2022). More recent algorithms, such as DPO (Rafailov et al., 2023) and CPO (Xu et al., 2024a), which replace reinforcement learning (RL) in PPO with *preference* optimization (PO), are reminiscent of the pairwise ranking optimization approach to tuning (Hopkins and May, 2011). Given this close relationship between these approaches, we can consider whether other techniques developed to improve MT tuning could be applied to optimization for style transfer.

In this work, we propose Style TrAnsfer with Multi-iteration Preference optimization (STAMP), a two-phase PO training framework, in which we

¹not to be confused with word alignment.

²not to be confused with parameter fine-tuning.

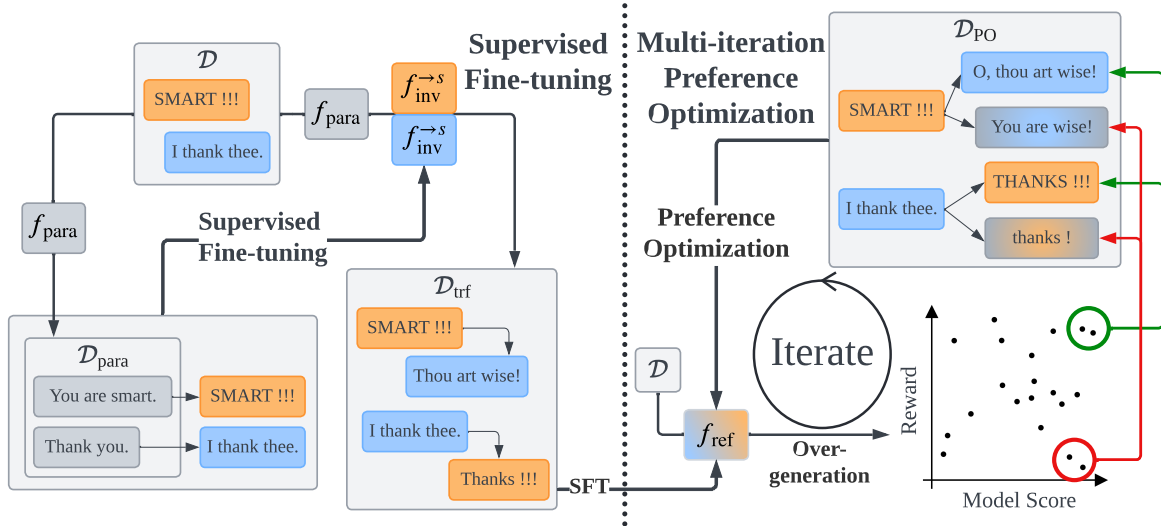


Figure 1: An overview of STAMP, in which we first train a unified style transfer model using supervised fine-tuning on pseudo-parallel data generated from non-parallel data, and then further train the model using multi-iteration preference optimization on preference pairs constructed with hope-and-fear sampling.

079 first use supervised fine-tuning to build a reference
 080 model from pseudo-parallel data and then train the
 081 reference model using PO. STAMP is similar to
 082 STEER and ASTRAPOP at a high level but is en-
 083 hanced with two techniques borrowed from MT
 084 tuning and two modifications that further adapt it
 085 for text style transfer. First, we include *multiple*
 086 *iterations* of preference pair generation followed
 087 by model optimization (Och, 2003), which has al-
 088 ready been shown to be effective on other Seq2Seq
 089 tasks such as mathematical and scientific reasoning
 090 (Chen et al., 2024; Pang et al., 2024; Song et al.,
 091 2024b; Yuan et al., 2024). Second, following the
 092 hope-and-fear sampling in Chiang (2012), for PO,
 093 we over-generate outputs using the reference model
 094 and construct preference pairs using samples with
 095 high model scores and extreme (high or low) task
 096 objective scores, in order to avoid dangerous gen-
 097 eration and encourage reachable good generation.
 098 To improve the quality of the reference model and
 099 the balance across the multiple training objectives,
 100 we additionally design a new two-step end-to-end
 101 pseudo-parallel data generation method and a dy-
 102 namic reward aggregation method.

103 We evaluate our model on two popular text style
 104 transfer datasets, Grammarly’s Yahoo Answers For-
 105 mality Corpus (GYAFC) (Rao and Tetreault, 2018)
 106 and the Corpus of Diverse Styles (CDS) (Krishna
 107 et al., 2020). Extensive experiments show that our
 108 model performs well on both in-domain and out-
 109 of-domain text style transfer, and outperforms all
 110 state-of-the-art baselines on both datasets.

111 Our main contributions are:

- We propose a multi-iteration contrastive preference optimization training framework with hope-and-fear preference pair construction for text style transfer.
- We design a new pseudo-parallel generation strategy and a dynamic weighted rewarded aggregation method to enhance the training framework for text style transfer.
- We show that, with the enhancements, our training framework produces style transfer models that achieve state-of-the-art performance on two popular text style transfer datasets.³

2 Methodology

In this section, we formalize the text style transfer task and introduce our training framework, STAMP.

2.1 Task Definition

Given a source text \mathbf{x} and a desired target style s , the goal of text style transfer is to generate a fluent rewrite of \mathbf{x} , denoted as $\mathbf{x} \rightarrow^s$, that has the same meaning as \mathbf{x} but is in style s . In this work, we focus on high-resource text style transfer in which we have access to a reasonable number of texts⁴ for each target style. Specifically, we have a set of texts with style labels, denoted as $\mathcal{D} = \{(\mathbf{x}_1, s_1), \dots, (\mathbf{x}_n, s_n)\}$, where \mathbf{x}_i and s_i re-

³We will release our code, models, and data to enable reproduction studies.

⁴In this work, we assume at least 2000 texts per style.

fer to the i^{th} text and its style, respectively. For convenience, we adopt notations from Hallinan et al. (2023a) and denote the **fluency** of a text \mathbf{x}_i as $F(\mathbf{x}_i)$, the **meaning similarity** between two texts \mathbf{x}_i and \mathbf{x}_j as $MS(\mathbf{x}_i, \mathbf{x}_j)$, and the **target style strength** of a text \mathbf{x}_i w.r.t. a target style s as $TSS(\mathbf{x}_i, s)$. Thus, given \mathcal{D} , we aim to build a text style transfer system that maximizes three independent objectives: $F(\mathbf{x}^{\rightarrow s})$, $MS(\mathbf{x}, \mathbf{x}^{\rightarrow s})$, and $TSS(\mathbf{x}^{\rightarrow s}, s)$.⁵

2.2 Framework Overview

STAMP is a preference optimization-based training framework that contains two main stages, a supervised fine-tuning (SFT) stage and a multi-iteration preference optimization (PO) stage. In the SFT stage, we first generate a dataset \mathcal{D}_{trf} of end-to-end pseudo-parallel style transfer pairs from the (non-parallel) dataset \mathcal{D} and then train a style transfer model f_{SFT} on \mathcal{D}_{trf} using supervised fine-tuning. In the PO stage, we train a model initialized to f_{SFT} using multi-iteration PO⁶ to directly maximize the three objectives, TSS, MS, and F, and obtain our final transfer model f_{PO} .

2.3 Supervised Fine-tuning

Due to a lack of parallel data, we adopt the technique described by Krishna et al. (2020), in which style-oriented paraphrasing is used to generate pseudo-parallel transfer data for each target style. Specifically, we paraphrase the texts in \mathcal{D} using a general paraphraser f_{para} similar to Krishna et al. (2020) and Hallinan et al. (2023a). To ensure meaning similarity preservation of the paraphrases, we generate k_{para} paraphrases for each text $\mathbf{x}_i \in \mathcal{D}$ and select the one with the highest meaning similarity to the original text, denoting it \mathbf{p}_i . We then obtain a dataset of paraphrases $\mathcal{D}_{\text{para}} = \{\mathbf{p}_1, \dots, \mathbf{p}_n\}$. For each target style s , we train a Seq2Seq model $f_{\text{inv}}^{\rightarrow s}$ ⁷ on $\{(\mathbf{p}_i \rightarrow \mathbf{x}_i) \mid 0 \leq i \leq n \text{ and } s_i = s\}$ to maximize

$$p(\mathbf{x} \mid \mathbf{p}) = \prod_{i=1}^{|\mathbf{x}|} p(\mathbf{x}[i] \mid \mathbf{p}, \mathbf{x}[\leq i]) \quad (1)$$

where $\mathbf{x}[i]$ and $\mathbf{x}[\leq i]$ represent the i^{th} token in \mathbf{x} and tokens preceding the i^{th} token in \mathbf{x} , respectively.

Following Krishna et al. (2020), we can transfer the style of a text \mathbf{x} to a style s through

$$\mathbf{x}^{\rightarrow s} = f_{\text{inv}}^{\rightarrow s}(f_{\text{para}}(\mathbf{x})) \quad (2)$$

⁵For brevity, we omit the arguments where unambiguous.

⁶See § 3.5 for details on the choice of PO used here.

⁷‘inverse’ due to data provenance, c.f. (Krishna et al., 2020)

where $\mathbf{x}^{\rightarrow s}$ is the transferred text. However, the two-step generation breaks the gradient connection between \mathbf{x} and $\mathbf{x}^{\rightarrow s}$ which is needed in the PO stage to maximize the meaning similarity between \mathbf{x} and $\mathbf{x}^{\rightarrow s}$. Therefore, we need an end-to-end pseudo-parallel dataset \mathcal{D}_{trf} to train a model that directly transfers a source text to each target style with no intermediate step.

To obtain \mathcal{D}_{trf} , we transfer the texts in \mathcal{D} using f_{para} and $f_{\text{inv}}^{\rightarrow s}$ for each target style s . Specifically, for each target style s , we transfer the texts in other styles in \mathcal{D} using Eq. 2 and obtain a dataset of style transfer pairs $\mathcal{D}_{\text{trf}}^{\rightarrow s} = \{(\mathbf{x}_i \rightarrow \mathbf{t}_i, s) \mid (\mathbf{x}_i, s_i) \in \mathcal{D} \text{ and } s_i \neq s\}$, where $\mathbf{t}_i = f_{\text{inv}}^{\rightarrow s}(f_{\text{para}}(\mathbf{x}_i))$ is a transfer of \mathbf{x}_i in style s . To obtain high-quality transferred texts, we generate k_{sft} transfers for each source text and select the one with the highest $F \cdot MS^{\tau_{\text{ms}}} \cdot TSS$, where $\tau_{\text{ms}} > 1$ is a temperature hyperparameter incorporated into the MS term to emphasize meaning similarity. We then construct \mathcal{D}_{trf} by combining $\mathcal{D}_{\text{trf}}^{\rightarrow s}$ for all target styles and train an end-to-end style transfer model f_{SFT} on the combined data \mathcal{D}_{trf} to maximize

$$p(\mathbf{t} \mid \mathbf{x}) = \prod_{i=1}^{|\mathbf{t}|} p(\mathbf{t}[i] \mid \mathbf{x}, \mathbf{t}[\leq i], s) \quad (3)$$

Note that unlike Eq. 2, the probability in Eq. 3 is also conditioned on s because we adopt the unified model setting in (Hallinan et al., 2023a). That is, we have a single transfer model for all target styles and control the target style with control codes.

2.4 Multi-iteration Preference Optimization

We further train the SFT model f_{SFT} from the previous stage with multi-iteration PO to directly optimize the model on the style transfer objectives: F, MS, and TSS. To apply PO (Rafailov et al., 2023; Xu et al., 2024a) we first generate paired preference data from a *reference model* f_{ref} and then train a model on this offline preference data in a contrastive manner starting from the reference model. Inspired by Och (2003) and recent studies in iterative PO, such as Yuan et al. (2024) and Chen et al. (2024), we perform PO for multiple iterations to improve over the offline-only training, updating the reference model between iterations. Specifically, in iteration i , we construct preference dataset $\mathcal{D}_{\text{PO}}^i$ by transferring texts drawn from \mathcal{D} , using reference model f_{ref}^i . We use PO (Rafailov et al., 2023; Xu et al., 2024a) to train a model initialized to f_{ref}^i to match the preferences in $\mathcal{D}_{\text{PO}}^i$; we

call the resulting model f_{PO}^i . We define f_{ref}^1 to be f_{SFT} and in all other cases we define f_{ref}^i to be f_{PO}^{i-1} . We next detail how the preference pairs in $\mathcal{D}_{\text{PO}}^i$ are constructed and the reward function used in this process.

2.4.1 PO Data Generation

We construct the preference dataset from \mathcal{D} using the hope-and-fear sampling strategy in Chiang (2012). While that work used BLEU (Papineni et al., 2002) as a preference metric, we instead use our style transfer reward \mathcal{R} which is detailed in § 2.4.2. Specifically, for each style s , we generate k_{PO} rewrites of each text \mathbf{x}_i in \mathcal{D} , whose initial style $s_i \neq s$, into style s and select the preference pair from the rewrites based on both the reward scores \mathcal{R} and the model scores \mathcal{M} of the rewrites, where \mathcal{M} is the average token-level probability w.r.t. f_{ref} . We select the rewrite with the highest $\mathcal{M}^{\tau_{\mathcal{M}}} + \mathcal{R}$ as the “winning” rewrite \mathbf{t}_i^w and the rewrite with the highest $\mathcal{M}^{\tau_{\mathcal{M}}} - \mathcal{R}$ as the “losing” rewrite⁸ \mathbf{t}_i^l , where $\tau_{\mathcal{M}}$ is the temperature controlling the weight of model score.⁹ We then obtain a new dataset $\mathcal{D}_{\text{PO}}^{\rightarrow s} = \{(\mathbf{x}_i \rightarrow (\mathbf{t}_i^w, \mathbf{t}_i^l), s) \mid (\mathbf{x}_i, s_i) \in \mathcal{D}\}$ for each style s . Combining $\mathcal{D}_{\text{PO}}^{\rightarrow s}$ for all styles, we finally obtain the PO dataset \mathcal{D}_{PO} .

2.4.2 Reward Function

To directly maximize the three objectives, F, MS, and, TSS, we use an aggregation of them as the reward function \mathcal{R} . The most straightforward aggregation is to take the product of the three as in Hallinan et al. (2023a). However, since the three objectives are independent, the probability of generating samples that have high scores in all three objectives is very low. Our preliminary experiments show that samples with high total rewards can also have low single-objective scores, which naturally results in preference pairs in which the “winning” outputs have lower single-objective scores. We refer to these as *reversed single-objective scores*. When the percentage of reversed single-objective scores is high, we observe a degradation in the corresponding objective after PO. To prevent the degradation in any objective, we propose to use a weighted product, which is given by

$$\mathcal{R} = \text{TSS}^\alpha \cdot \text{MS}^\beta \cdot \text{F}^\gamma \quad (4)$$

⁸also called “chosen” and “rejected” rewrites in PO literature (e.g., Rafailov et al., 2023).

⁹In practice, we find using model score does not benefit performance, so we drop this term for STAMP, which reduces the preference pair selection criteria to the sample with the highest \mathcal{R} and $-\mathcal{R}$; a detailed comparison is shown in § 4.3.

where α , β , and γ are temperatures for each objective.

We dynamically calculate α , β , and γ based on the number of reversed single-objective scores in the preference pairs for each iteration. For convenience, we denote the number of reversed single-objective scores for each objective as r_{TSS} , r_{MS} , and r_{F} .¹⁰ We first set $\beta = \gamma = 1$ and set α to be the smallest positive integer such that $r_{\text{TSS}} < r_{\text{MS}}$ and $r_{\text{TSS}} < r_{\text{F}}$. Then, we fix α and γ and set β to be the largest positive integer such that $r_{\text{MS}} > r_{\text{TSS}}$. Finally, we fix α and β and set γ to be the largest positive integer such that $r_{\text{F}} > r_{\text{TSS}}$ and $r_{\text{F}} > r_{\text{MS}}$. We set an upper bound τ_{max} to α , β , and γ to prevent \mathcal{R} from leaning too much to any objective.

3 Experiments

We evaluate STAMP on two text style transfer datasets in both in-domain and out-of-domain settings and compare STAMP with the state-of-the-art baseline approaches. In this section, we detail the experimental setup and the model implementation.

3.1 Datasets

We use two style transfer datasets in this work: (1) **Corpus of Diverse Styles (CDS)** (Krishna et al., 2020), which contains non-parallel texts in 11 different styles, such as Shakespeare and English Tweets, and (2) **Grammarly’s Yahoo Answers Formality Corpus (GYAFC)** (Rao and Tetreault, 2018), which contains non-parallel formal and informal texts for training and a small number of parallel transfer pairs for tuning and test. In this work, we only use non-parallel texts with style labels for training, validation, and test.

To reduce computational costs, we use a subset of each dataset. Specifically, we sample 2000 texts per style for training, and 200 per style for validation. For CDS we sample 200 per style for test, while for GYAFC we sample 1000 per style. When constructing the end-to-end pseudo-parallel dataset \mathcal{D}_{trf} , for each target style, we sample 200 and 20 source texts from each of the other styles for training and validation, respectively. In the in-domain testing, we transfer the test texts in each style to all other styles in the same dataset and calculate the total average scores and average scores grouped by the target style. In the out-of-domain testing, we transfer all test texts in each dataset to all styles in

¹⁰ r_{TSS} , r_{MS} , and r_{F} are functions of α , β , and γ , so we recalculate r ’s each time we change the value of α , β , or γ .

the other dataset and calculate the same scores. We elaborate on metric scores in § 3.4.1.

Besides the style transfer datasets, we also use a paraphrase dataset, **ParaNMT** (Wieting and Gimpel, 2018) to train the paraphraser used for pseudo-parallel data generation. Specifically, we use the filtered version containing 75k paraphrase pairs in Krishna et al. (2020).

3.2 Reward Models

We have a reward model for each of the three objectives, TSS, MS, and F. For convenience, we use the same notations to refer to the objective functions and the corresponding reward models in this paper.

Target Style Strength (TSS) We use a single style classifier, f_{cls} with multiple binary sigmoid classification heads to calculate the TSS for each target style. We train f_{cls} from the pre-trained RoBERTa-large model (Liu et al., 2019b) on the same training and validation splits as discussed in § 3.1. We simply use the sigmoid outputs from the classification heads as the TSS scores which range from 0 to 1.

Meaning Similarity (MS) We assess the meaning similarity between the source text and the transferred text using the cosine similarity between the semantic embeddings of the two texts. The semantic embeddings are calculated using SBERT¹¹ (Reimers and Gurevych, 2019). Technically, the cosine similarity of two embeddings ranges from -1 to 1, but negative cosine similarity is very rare in our experiments since we always the similarity between two paraphrases. Following Hallinan et al. (2023a), we clip negative values to 0 to ensure that MS ranges from 0 to 1.

Fluency (F) To measure the fluency of a text, we use a text classifier¹² trained on the Corpus of Linguistic Acceptability (CoLA) (Warstadt et al., 2019). The softmax score of the “grammatical” class is used as the F score which also ranges from 0 to 1.

3.3 Baseline Approaches

We compare STAMP with 4 strong baselines: GPT prompting (Reif et al., 2022), STRAP (Krishna et al., 2020), STEER (Hallinan et al., 2023a), and ASTRAPOP (Liu et al., 2024).

¹¹We use the variant with the best sentence embedding performance, which is all-mpnet-base-v2.

¹²<https://huggingface.co/cointegrated/roberta-large-cola-krishna2020>

GPT prompting uses the zero- and few-shot capability of GPT-3.5-turbo to transfer texts to the target style given just the name of the style and 5 target style exemplars (5-shot) or no exemplars (zero-shot).

STRAP transfers a text by paraphrasing the text with a diverse paraphraser followed by an inverse paraphraser trained on pseudo-parallel transfer data generated by the diverse paraphraser.

STEER generates pseudo-parallel data using an expert-guided generation technique (Liu et al., 2021a), and trains an end-to-end style transfer model on the generated data using a reinforcement learning algorithm (Lu et al., 2022).

ASTRAPOP adopts the same paraphrase-and-inverse-paraphrase pipeline as STRAP but trains the inverse paraphraser using policy optimization or PO to directly maximize the target style strength, which achieves better performance on both low-resource and high-resource authorship style transfer. It does not use multi-iteration optimization, nor the overgeneration strategies we describe.

3.4 Evaluation Metrics

3.4.1 Automatic Evaluation

We evaluate the approaches on the three objectives, TSS, MS, and F, using the same reward models introduced in § 3.2. To assess overall performance, we use a single aggregate score $\text{Agg.} = \text{TSS} \cdot \text{MS} \cdot \text{F}$. Note that the reward models described in § 3.2 calculate scores for single transfer pairs, while the final scores used for evaluation are averages over all transfer pairs in the test set.

3.4.2 Human Evaluation

In addition to the automatic evaluation, we conduct a human evaluation to assess the model performance on the three style transfer objectives: TSS_h , MS_h , and F_h .¹³ For TSS_h , we show 5 exemplars for the style of the input text and 5 exemplars for the target style, and ask the annotator to select the style of the transferred text out of these two styles. The sample gets a score of 1 if the target style is selected, and 0 otherwise. For MS_h and F_h , we ask whether the transferred text has a similar meaning to the input text and whether the transferred is fluent, respectively, and collect the answers using a three-level Likert scale ranging from 0 to 2. See § B.5 for the detailed instructions used in the human evaluation.

¹³We use the subscript h to distinguish human metrics from automatic metrics.

Approach	CDS				GYFAC			
	TSS	MS	F	Agg.	TSS	MS	F	Agg.
GPT zero-shot	0.189 [‡]	0.705 [‡]	0.803 [†]	0.104 [‡]	0.672 [‡]	0.788 [‡]	0.968	0.489 [‡]
GPT 5-shot	0.199 [‡]	<u>0.735</u> [†]	<u>0.805</u> [†]	0.112 [‡]	0.667 [‡]	<u>0.800</u> [†]	<u>0.965</u>	0.495 [‡]
STRAP	0.382 [‡]	0.626 [‡]	0.759 [‡]	0.158 [‡]	0.618 [‡]	0.735 [‡]	0.913 [‡]	0.409 [‡]
STEER	<u>0.654</u> [†]	0.672 [‡]	0.905	<u>0.395</u> [†]	0.951	0.776 [‡]	0.930 [‡]	<u>0.686</u> [†]
ASTRAPOP	0.542 [‡]	0.600 [‡]	0.755 [‡]	0.221 [‡]	0.783 [‡]	0.734 [‡]	0.924 [‡]	<u>0.525</u> [‡]
STAMP	0.746	0.801	0.801 [†]	0.474	0.958	0.921	0.941 [‡]	0.828

Table 1: The automatic evaluation results on in-domain inputs on the CDS and the GYFAC datasets. The best and the 2nd best scores in each column are shown in **bold** and underline, respectively. “†” and “‡” indicate the score is significantly ($p < 0.05$) worse than the best score and the top 2 scores in the same column, respectively, determined by resampling t-test.

3.5 Implementation Details

We implement all Seq2Seq models in STAMP, including the paraphraser and all transfer models, as decoder-only Seq2Seq models (Wolf et al., 2019) based on pre-trained LLaMA-2-7B (Touvron et al., 2023). The input and output are concatenated together with a separator token “[SEP].” For the unified transfer model f_{SFT} , we prepend a style code for the target style (e.g., “[SHAKESPEARE]” and “[FORMAL]”) to the input to control the output style. We use CPO (Xu et al., 2024a) in the multi-iteration PO stage. We choose CPO instead of the most popular PO algorithm, DPO (Rafailov et al., 2023), since CPO has been shown to be more efficient and effective (Xu et al., 2024a; Liu et al., 2024). Also, compared to DPO, CPO has an additional negative log-likelihood term that is found to be significant for multi-iteration preference optimization (Pang et al., 2024). We stop PO training at the iteration where the validation TSS starts to decrease and use the model from the previous iteration as the final model. For fairness, all non-GPT baselines are also implemented based on LLaMA-2-7B and use the same paraphraser as STAMP. We use gpt-3.5-turbo-0125 for all GPT-based approaches. See § B for hyperparameters and GPT zero- and few-shot prompts.

4 Results

In this section, we present the quantitative experimental results. A qualitative case study is in § A.3. Because of the limited resources, we conduct all experiments for a single run and perform t-tests on the results.¹⁴

¹⁴See § B.1 for details.

4.1 Automatic Evaluation

Automatic evaluation results on in-domain input are shown in Table 1. According to the aggregated score (Agg.), STAMP outperforms all baselines on the overall performance by a large margin on both datasets. Looking at the per-objective scores, STAMP has the best target style strength (TSS) and meaning similarity (MS), but its fluency (F) is relatively lower, and this disadvantage is more obvious on the CDS dataset. STEER has the best overall performance (Agg.) among the baselines on both datasets, while the overall performance of other baselines are mixed across the two datasets. The results on the out-of-domain style transfer experiments are generally consistent with the in-domain results. See § A.1 for details.

Approach	TSS	MS _h /2	F _h /2	Agg. _{~h}
GPT 5-shot	0.16	<u>0.75</u>	<u>0.90</u>	0.11
STEER	<u>0.58</u>	<u>0.62</u>	0.92	0.33
STAMP	0.79	0.75	0.80	0.47

Table 2: The human evaluation results on in-domain inputs on the CDS datasets. The best and the 2nd best scores in each column are shown in **bold** and underline, respectively.

4.2 Human Evaluation

We conduct a human evaluation on the CDS dataset for STAMP, the best-performing baseline (STEER), and the best GPT-prompting baseline (GPT 5-shot). We randomly choose 5 samples from each of the 11 target styles for each of the three models, which yields 165 samples in total, and collect up to three annotations for each sample. Seven volunteer NLP experts are recruited for annotation. We perform an independent sample t-test on the annotation results and find statistically significant differences

Approach	CDS				GYFAC			
	TSS	MS	F	Agg.	TSS	MS	F	Agg.
STAMP	0.746	0.801 [‡]	<u>0.801</u> [†]	0.474	0.958 [‡]	0.921 [†]	0.941 [†]	0.828
$\tau_{\mathcal{M}} = 0.1$	0.720 [†]	0.796 [‡]	0.800 [†]	0.454 [†]	0.965	0.910 [‡]	0.943 [†]	0.826
$k_{\text{PO}} = 2$	0.745	0.688 [‡]	0.816	0.411 [‡]	0.970	0.878 [‡]	0.947	0.804 [‡]
Random t^l	0.640 [‡]	0.836	0.780 [‡]	0.412 [‡]	0.950 [‡]	0.924 [†]	0.937	0.822
High t^l	0.592 [‡]	<u>0.826</u> [†]	0.796 [†]	0.384 [‡]	0.928 [‡]	0.936	0.932 [‡]	0.810 [‡]

Table 3: Hope-and-fear sampling ablations, evaluated automatically on in-domain inputs on the CDS and the GYFAC datasets. The best and the 2nd best scores in each column are shown in **bold** and underline, respectively. “†” and “‡” indicate the score is significantly ($p < 0.05$) worse than the best score and the top 2 scores in the same column, respectively, determined by resampling t-test.

in MS_h and F_h but not in TSS_h ,¹⁵ which is in line with our expectation since the style classification has been found to be hard for untrained humans (Krishna et al., 2020; Hallinan et al., 2023a). Therefore, following Krishna et al. (2020) and Hallinan et al. (2023a), we calculate the quasi aggregated score $\text{Agg.}_{\sim h}$ using TSS, MS_h , and F_h . Formally, $\text{Agg.}_{\sim h} = \text{TSS} \cdot \frac{MS_h}{2} \cdot \frac{F_h}{2}$, where we divide MS_h and F_h by 2 to scale them to the $[0, 1]$ range so that $\text{Agg.}_{\sim h}$ also ranges from 0 to 1. As shown in Table 2, STAMP has the best meaning similarity (MS_h) and overall performance ($\text{Agg.}_{\sim h}$), but its fluency is worse than STEER and GPT 5-shot transfer, which is consistent with the automatic evaluation results.

4.3 Ablation Studies

In this section, we demonstrate the effects of our four main contributions in STAMP: multi-iteration PO, hope-and-fear sampling, weighted reward aggregation, and end-to-end pseudo-parallel data generation.

Multi-iteration PO & Weighted \mathcal{R} We show the performance evolution of STAMP and STAMP with unweighted \mathcal{R} over the multi-iteration PO training in Figure 2. In general, the overall performance (Agg.) of both models keeps increasing over the iterations, which indicates the effectiveness of multi-iteration optimization. STAMP with unweighted \mathcal{R} performs slightly better than STAMP, but it has a severe degradation in meaning similarity (MS), and the scores in the three objectives have a substantial difference after training. In contrast, with the weighted reward aggregation, STAMP shows a higher stability in all scores. Only

¹⁵See § A.2 for the raw human evaluation scores and the result of the t-test.

¹⁶which is calculated from the human study samples using the automatic TSS metric.

fluency (F) exhibits a slight decrease, and scores in all three objectives converge to a similar value at the end of the training.

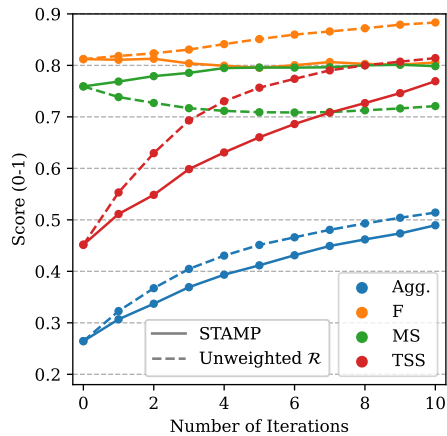


Figure 2: The value of iterative CPO on performance in STAMP and STAMP with unweighted \mathcal{R} , shown on the CDS dataset (test split). Iteration 0 refers to the SFT model before PO.

Hope-and-fear Sampling The results of hope-and-fear sampling ablation are shown in Table 3. As mentioned in § 2.4.2, we do not use the model score term in hope-and-fear sampling for preference pair construction since it does not improve the performance, which can be observed from the “ $\tau_{\mathcal{M}} = 0.1$ ” row in Table 3. The last three rows in Table 3 show that both dropping over-generation ($k_{\text{PO}} = 2$) and using a random other sample (Random t^l) or the sample with the second highest reward (High t^l) as the “losing” sample undermine the overall performance of STAMP.

Pseudo-parallel Data Generation We demonstrate the superiority of our two-step end-to-end pseudo-parallel data generation method by comparing the STAMP SFT model, f_{SFT} , with the best-performing baseline SFT style transfer model, STRAP. The overall performance (Agg.) of the two

models is shown in Table 4. With our method, the overall performance of f_{SFT} is much higher than STRAP on both datasets, which provides a better starting point for PO.

	CDS	GYAFC
STRAP	0.158	0.409
f_{SFT}	0.264	0.657

Table 4: The overall performance (Agg.) of STRAP and the STAMP SFT model (f_{SFT}) on CDS and GYAFC. The best score in each column is shown in **bold**.

5 Related Work

Text Style Transfer Due to the lack of parallel style transfer data, only a limited number of studies address this task as a supervised or semi-supervised Seq2Seq task, which requires a certain amount of parallel data for training and/or tuning (Zhu et al., 2010; Rao and Tetreault, 2018; Wang et al., 2019; Shang et al., 2019; Xu et al., 2019; Zhang et al., 2020; Kim et al., 2022; Raheja et al., 2023). Although these approaches work well when parallel data is available, none generalize well to styles with no parallel data. As a result, most works in this area focus on unsupervised approaches that require only non-parallel data or even no data. These works mainly approach the task via latent representation disentanglement and manipulation (Lample et al., 2019; Liu et al., 2019a; John et al., 2019; Jin et al., 2020), style-related pattern editing (Madaan et al., 2020; Malmi et al., 2020; Reid and Zhong, 2021; Luo et al., 2023), pseudo-parallel transfer data construction (Krishna et al., 2020; Riley et al., 2021), policy optimization (Gong et al., 2019; Liu et al., 2021b; Deng et al., 2022; Hallinan et al., 2023a; Liu et al., 2024), and LLM zero- or few-shot prompting (Reif et al., 2022; Suzgun et al., 2022; Patel et al., 2023).

Among these approaches, two of the policy optimization based approaches, STEER (Hallinan et al., 2023a) and ASTRAPOP (Liu et al., 2024) achieve the best performance on text style transfer and authorship style transfer, respectively. Their high-level training frameworks both combine pseudo-parallel data generation and policy optimization, but their specific approaches differ. For pseudo-parallel data generation, STEER uses a paraphraser guided by an expert and an anti-expert, while ASTRAPOP simply paraphrases the texts in the target style and uses these paraphrase-to-target transfer pairs. For policy optimization, STEER uses an RL

algorithm, Quark, while ASTRAPOP tries three options: one RL algorithm, PPO (Schulman et al., 2017), and two PO algorithms, DPO (Rafailov et al., 2023) and CPO (Xu et al., 2024a). Our framework shares the same high-level procedure with STEER and ASTRAPOP, but we design a new pseudo-parallel data generation method and also enhance the PO stage with multi-iteration training, weighted reward aggregation, and hope-and-fear preference pair construction. These enhancements dramatically improve the performance of STAMP over STEER and ASTRAPOP.

Preference Optimization PO (Rafailov et al., 2023; Song et al., 2024a; Xu et al., 2024a) is a class of RL-free policy optimization algorithms which has been broadly applied to train generative language models on direct task objectives instead of the language modeling loss and is closely related to (pre-neural) machine translation objective ‘tuning’ (Och, 2003; Chiang et al., 2008; Hopkins and May, 2011). Rafailov et al. (2023) show that PO is more stable and efficient than traditional RL-based algorithms on sentiment generation and text summarization (Rafailov et al., 2023). It has also been successfully applied to many other NLP tasks, such as training helpful and harmless assistants (Song et al., 2024a), machine translation (Xu et al., 2024a), and authorship style transfer (Liu et al., 2024). Later works (Xiong et al., 2023; Xu et al., 2024b; Yuan et al., 2024; Chen et al., 2024; Pang et al., 2024; Song et al., 2024b) extend the offline PO algorithms by performing the optimization for multiple iterations and further improve the performance of the models. In this work, we adopt the multi-iteration PO for STAMP and enhance it with weighted reward aggregation and hope-and-fear preference pair construction, which improve the effectiveness of multi-iteration PO training.

6 Conclusion

We present STAMP, a multi-iteration preference optimization training framework for text style transfer, in which an end-to-end pseudo-parallel data generation pipeline provides a strong reference model, a preference pair construction strategy improves the effectiveness of PO training, and weighted reward aggregation ensures balance across multiple objectives over multi-iteration training. We evaluate STAMP on two commonly used text style transfer datasets; demonstrating superior performance over all state-of-the-art style transfer approaches.

624 Limitations

625 Although achieving the state-of-the-art perfor-
626 mance on two text style transfer datasets, STAMP
627 has two main limitations. First, we observe rep-
628 etitions and hallucinations in some transferred
629 texts. The potential reason is that PO training in-
630 creases the peakiness of the model, which means
631 the probability of generating the tokens that are
632 frequent in the target style increases dispropor-
633 tionately (Choshen et al., 2020; Kiegedland and
634 Kreutzer, 2021). The occurrence of repetitions and
635 hallucinations also indicates that our reward model
636 cannot fully capture all aspects of the desired ob-
637 jectives. Two possible solutions are developing PO
638 algorithms that are less vulnerable to the increased
639 peakiness and developing better reward models.
640 These are two promising directions for future stud-
641 ies but are out of the scope of the current work
642 which focuses on the multi-iteration extension of
643 existing preference optimization algorithms and the
644 strategies for preference pair construction.

645 Second, as discussed in § 4.3, the weighted re-
646 ward aggregation method is effective on the CDS
647 dataset but is not very useful on the GYAFC dataset
648 because formality transfer is a relatively easier task,
649 and it is more likely to generate high-quality sam-
650 ples with balanced single-objective scores. It could
651 be useful to add a control mechanism to determine
652 when using the weighted aggregation is beneficial
653 to prevent overbalanced single-objective scores on
654 easy tasks.

655 Ethical Considerations

656 As a general text style transfer framework, STAMP
657 can transfer texts to any target style given an ade-
658 quate amount of non-parallel data, which means it
659 can potentially be used to generate unethical texts
660 such as transferring normal texts into an offensive
661 or profane style. Moreover, although STAMP is
662 not specifically designed for authorship transfer, it
663 can still serve that purpose by transferring the texts
664 into the style of a particular author, which can be
665 unethical if used without authorization. However,
666 privatization of an author’s style can also be used
667 to enable oppressed people to communicate freely
668 without the fear of recrimination. In any case, as
669 we and others show, the state of the art of style
670 transfer is not yet advanced for either privacy or
671 mimicry to be a significant concern in a deployed
672 system. Our work is strictly intended for research
673 and personal use on public or authorized data.

674 Some texts in the datasets used in this work
675 (though collected and released elsewhere) contain
676 words or ideas that may cause harm to others. We
677 do not generally filter out those texts, so that we
678 may maximally preserve the characteristics of dif-
679 ferent styles. However, for human studies, we
680 remove all texts with personal identifiable infor-
681 mation (PII) to ensure privacy and remove texts
682 that contain profane language to minimize harm
683 to human subjects. We exclude these texts in-
684 stead of masking out PII or profane tokens, since
685 masks may influence annotators’ judgments regard-
686 ing meaning similarity and fluency. The protocols
687 of our human studies have been approved by an
688 institutional review board.

References 689

- 690 Haohan Bo, Steven H. H. Ding, Benjamin C. M. Fung,
691 and Farkhund Iqbal. 2021. [ER-AE: Differentially
692 private text generation for authorship anonymization.](#)
693 In *Proceedings of the 2021 Conference of the North
694 American Chapter of the Association for Computa-
695 tional Linguistics: Human Language Technologies*,
696 pages 3997–4007, Online. Association for Computa-
697 tional Linguistics.
- 698 Zixiang Chen, Yihe Deng, Huizhuo Yuan, Kaixuan Ji,
699 and Quanquan Gu. 2024. [Self-play fine-tuning con-
700 verts weak language models to strong language mod-
701 els.](#) *Preprint*, arXiv:2401.01335.
- 702 Colin Cherry and George Foster. 2012. [Batch tuning
703 strategies for statistical machine translation.](#) In *Pro-
704 ceedings of the 2012 Conference of the North Amer-
705 ican Chapter of the Association for Computational
706 Linguistics: Human Language Technologies*, pages
707 427–436, Montréal, Canada. Association for Compu-
708 tational Linguistics.
- 709 David Chiang. 2012. [Hope and fear for discriminative
710 training of statistical translation models.](#) *Journal of
711 Machine Learning Research*, 13(40):1159–1187.
- 712 David Chiang, Yuval Marton, and Philip Resnik. 2008.
713 [Online large-margin training of syntactic and struc-
714 tural translation features.](#) In *Proceedings of the 2008
715 Conference on Empirical Methods in Natural Lan-
716 guage Processing*, pages 224–233, Honolulu, Hawaii.
717 Association for Computational Linguistics.
- 718 Leshem Choshen, Lior Fox, Zohar Aizenbud, and Omri
719 Abend. 2020. [On the weaknesses of reinforcement
720 learning for neural machine translation.](#) In *Interna-
721 tional Conference on Learning Representations*.
- 722 David Dale, Anton Voronov, Daryna Dementieva, Var-
723 vara Logacheva, Olga Kozlova, Nikita Semenov, and
724 Alexander Panchenko. 2021. [Text detoxification us-
725 ing large pre-trained neural models.](#) In *Proceedings*

841	Yixin Liu, Graham Neubig, and John Wieting. 2021b.	Rafael Rafailov, Archit Sharma, Eric Mitchell, Christo-	898
842	On learning text style transfer with direct rewards.	phers D Manning, Stefano Ermon, and Chelsea Finn.	899
843	In <i>Proceedings of the 2021 Conference of the North</i>	2023. Direct preference optimization: Your language	900
844	<i>American Chapter of the Association for Computa-</i>	model is secretly a reward model. In <i>Thirty-seventh</i>	901
845	<i>tional Linguistics: Human Language Technologies,</i>	<i>Conference on Neural Information Processing Sys-</i>	902
846	pages 4262–4273, Online. Association for Computa-	<i>tems.</i>	903
847	tional Linguistics.		
848	Ximing Lu, Sean Welleck, Jack Hessel, Liwei Jiang,	Vipul Raheja, Dhruv Kumar, Ryan Koo, and Dongyeop	904
849	Lianhui Qin, Peter West, Prithviraj Ammanabrolu,	Kang. 2023. CoEdIT: Text editing by task-specific	905
850	and Yejin Choi. 2022. QUARK: Controllable text	instruction tuning. In <i>Findings of the Association</i>	906
851	generation with reinforced unlearning. In <i>Advances</i>	<i>for Computational Linguistics: EMNLP 2023,</i> pages	907
852	<i>in Neural Information Processing Systems.</i>	5274–5291, Singapore. Association for Computa-	908
		tional Linguistics.	909
853	Guoqing Luo, Yu Han, Lili Mou, and Mauajama Firdaus.	Sudha Rao and Joel Tetreault. 2018. Dear sir or madam,	910
854	2023. Prompt-based editing for text style transfer. In	may I introduce the GY AFC dataset: Corpus, bench-	911
855	<i>Findings of the Association for Computational Lin-</i>	marks and metrics for formality style transfer. In	912
856	<i>guistics: EMNLP 2023,</i> pages 5740–5750, Singapore.	<i>Proceedings of the 2018 Conference of the North</i>	913
857	Association for Computational Linguistics.	<i>American Chapter of the Association for Computa-</i>	914
858	Aman Madaan, Amrith Setlur, Tanmay Parekh, Barn-	<i>tional Linguistics: Human Language Technologies,</i>	915
859	abas Poczos, Graham Neubig, Yiming Yang, Ruslan	<i>Volume 1 (Long Papers),</i> pages 129–140, New Or-	916
860	Salakhutdinov, Alan W Black, and Shrimai Prabhu-	leans, Louisiana. Association for Computational Lin-	917
861	moye. 2020. Politeness transfer: A tag and generate	guistics.	918
862	approach. In <i>Proceedings of the 58th Annual Meet-</i>		
863	<i>ing of the Association for Computational Linguistics,</i>	Machel Reid and Victor Zhong. 2021. LEWIS: Lev-	919
864	pages 1869–1881, Online. Association for Computa-	enshtein editing for unsupervised text style transfer.	920
865	tional Linguistics.	In <i>Findings of the Association for Computational</i>	921
866	Eric Malmi, Aliaksei Severyn, and Sascha Rothe. 2020.	<i>Linguistics: ACL-IJCNLP 2021,</i> pages 3932–3944,	922
867	Unsupervised text style transfer with padded masked	Online. Association for Computational Linguistics.	923
868	language models. In <i>Proceedings of the 2020 Con-</i>		
869	<i>ference on Empirical Methods in Natural Language</i>	Emily Reif, Daphne Ippolito, Ann Yuan, Andy Coenen,	924
870	<i>Processing (EMNLP),</i> pages 8671–8680, Online. As-	Chris Callison-Burch, and Jason Wei. 2022. A recipe	925
871	sociation for Computational Linguistics.	for arbitrary text style transfer with large language	926
872	Franz Josef Och. 2003. Minimum error rate training	models. In <i>Proceedings of the 60th Annual Meet-</i>	927
873	in statistical machine translation. In <i>Proceedings</i>	<i>ing of the Association for Computational Linguistics</i>	928
874	<i>of the 41st Annual Meeting of the Association for</i>	<i>(Volume 2: Short Papers),</i> pages 837–848, Dublin,	929
875	<i>Computational Linguistics,</i> pages 160–167, Sapporo,	Ireland. Association for Computational Linguistics.	930
876	Japan. Association for Computational Linguistics.		
877	Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida,	Nils Reimers and Iryna Gurevych. 2019. Sentence-	931
878	Carroll Wainwright, Pamela Mishkin, Chong Zhang,	BERT: Sentence embeddings using Siamese BERT-	932
879	Sandhini Agarwal, Katarina Slama, Alex Ray, et al.	networks. In <i>Proceedings of the 2019 Conference on</i>	933
880	2022. Training language models to follow instruc-	<i>Empirical Methods in Natural Language Processing</i>	934
881	tions with human feedback. <i>Advances in neural in-</i>	<i>and the 9th International Joint Conference on Natu-</i>	935
882	<i>formation processing systems,</i> 35:27730–27744.	<i>ral Language Processing (EMNLP-IJCNLP),</i> pages	936
883	Richard Yuanzhe Pang, Weizhe Yuan, Kyunghyun Cho,	3982–3992, Hong Kong, China. Association for Com-	937
884	He He, Sainbayar Sukhbaatar, and Jason Weston.	putational Linguistics.	938
885	2024. Iterative reasoning preference optimization.		
886	<i>Preprint,</i> arXiv:2404.19733.	Parker Riley, Noah Constant, Mandy Guo, Girish	939
887	Kishore Papineni, Salim Roukos, Todd Ward, and Wei-	Kumar, David Uthus, and Zarana Parekh. 2021.	940
888	Jing Zhu. 2002. Bleu: a method for automatic evalu-	TextSETTR: Few-shot text style extraction and tun-	941
889	ation of machine translation. In <i>Proceedings of the</i>	able targeted restyling. In <i>Proceedings of the 59th</i>	942
890	<i>40th Annual Meeting of the Association for Computa-</i>	<i>Annual Meeting of the Association for Computational</i>	943
891	<i>tional Linguistics,</i> pages 311–318, Philadelphia,	<i>Linguistics and the 11th International Joint Confer-</i>	944
892	Pennsylvania, USA. Association for Computational	<i>ence on Natural Language Processing (Volume 1:</i>	945
893	Linguistics.	<i>Long Papers),</i> pages 3786–3800, Online. Association	946
894	Ajay Patel, Nicholas Andrews, and Chris Callison-	for Computational Linguistics.	947
895	Burch. 2023. Low-resource authorship style trans-		
896	fer: Can non-famous authors be imitated? <i>Preprint,</i>	John Schulman, Filip Wolski, Prafulla Dhariwal,	948
897	arXiv:2212.08986.	Alec Radford, and Oleg Klimov. 2017. Prox-	949
		imal policy optimization algorithms. <i>Preprint,</i>	950
		arXiv:1707.06347.	951
		Mingyue Shang, Piji Li, Zhenxin Fu, Lidong Bing,	952
		Dongyan Zhao, Shuming Shi, and Rui Yan. 2019.	953
		Semi-supervised text style transfer: Cross projection	954

955	in latent space. In <i>Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)</i> , pages 4937–4946, Hong Kong, China. Association for Computational Linguistics.	
956		
957		
958		
959		
960		
961	Rakshith Shetty, Bernt Schiele, and Mario Fritz. 2018. A4NT: Author attribute anonymity by adversarial training of neural machine translation . In <i>27th USENIX Security Symposium (USENIX Security 18)</i> , pages 1633–1650, Baltimore, MD. USENIX Association.	
962		
963		
964		
965		
966		
967	Feifan Song, Bowen Yu, Minghao Li, Haiyang Yu, Fei Huang, Yongbin Li, and Houfeng Wang. 2024a. Preference ranking optimization for human alignment. In <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> , volume 38, pages 18990–18998.	
968		
969		
970		
971		
972	Yifan Song, Da Yin, Xiang Yue, Jie Huang, Sujian Li, and Bill Yuchen Lin. 2024b. Trial and error: Exploration-based trajectory optimization for llm agents . <i>Preprint</i> , arXiv:2403.02502.	
973		
974		
975		
976	Mirac Suzgun, Luke Melas-Kyriazi, and Dan Jurafsky. 2022. Prompt-and-rerank: A method for zero-shot and few-shot arbitrary textual style transfer with small language models . In <i>Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing</i> , pages 2195–2222, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.	
977		
978		
979		
980		
981		
982		
983		
984	Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony Hartshorn, Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel Kloumann, Artem Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov, Pushkar Mishra, Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan Saladi, Alan Schelten, Ruan Silva, Eric Michael Smith, Ranjan Subramanian, Xiaoqing Ellen Tan, Binh Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zheng Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurelien Rodriguez, Robert Stojnic, Sergey Edunov, and Thomas Scialom. 2023. Llama 2: Open foundation and finetuned chat models . <i>Preprint</i> , arXiv:2307.09288.	
985		
986		
987		
988		
989		
990		
991		
992		
993		
994		
995		
996		
997		
998		
999		
1000		
1001		
1002		
1003		
1004		
1005		
1006		
1007	Yunli Wang, Yu Wu, Lili Mou, Zhoujun Li, and Wenhan Chao. 2019. Harnessing pre-trained neural networks with rules for formality style transfer . In <i>Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)</i> , pages 3573–3578, Hong Kong, China. Association for Computational Linguistics.	
1008		
1009		
1010		
1011		
1012		
1013		
1014		
	Alex Warstadt, Amanpreet Singh, and Samuel R. Bowman. 2019. Neural network acceptability judgments . <i>Transactions of the Association for Computational Linguistics</i> , 7:625–641.	1015
		1016
		1017
		1018
	John Wieting and Kevin Gimpel. 2018. ParaNMT-50M: Pushing the limits of paraphrastic sentence embeddings with millions of machine translations . In <i>Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 451–462, Melbourne, Australia. Association for Computational Linguistics.	1019
		1020
		1021
		1022
		1023
		1024
		1025
	Thomas Wolf, Victor Sanh, Julien Chaumond, and Clement Delangue. 2019. Transfertransfo: A transfer learning approach for neural network based conversational agents . <i>Preprint</i> , arXiv:1901.08149.	1026
		1027
		1028
		1029
	Wei Xiong, Hanze Dong, Chenlu Ye, Ziqi Wang, Han Zhong, Heng Ji, Nan Jiang, and Tong Zhang. 2023. Iterative preference learning from human feedback: Bridging theory and practice for rlhf under kl-constraint. In <i>ICLR 2024 Workshop on Mathematical and Empirical Understanding of Foundation Models</i> .	1030
		1031
		1032
		1033
		1034
		1035
		1036
	Haoran Xu, Amr Sharaf, Yunmo Chen, Weiting Tan, Lingfeng Shen, Benjamin Van Durme, Kenton Murray, and Young Jin Kim. 2024a. Contrastive preference optimization: Pushing the boundaries of llm performance in machine translation . <i>Preprint</i> , arXiv:2401.08417.	1037
		1038
		1039
		1040
		1041
		1042
	Jing Xu, Andrew Lee, Sainbayar Sukhbaatar, and Jason Weston. 2024b. Some things are more cringe than others: Iterative preference optimization with the pairwise cringe loss . <i>Preprint</i> , arXiv:2312.16682.	1043
		1044
		1045
		1046
	Ruochen Xu, Tao Ge, and Furu Wei. 2019. Formality style transfer with hybrid textual annotations . <i>Preprint</i> , arXiv:1903.06353.	1047
		1048
		1049
	Weizhe Yuan, Richard Yuanzhe Pang, Kyunghyun Cho, Xian Li, Sainbayar Sukhbaatar, Jing Xu, and Jason Weston. 2024. Self-rewarding language models . <i>Preprint</i> , arXiv:2401.10020.	1050
		1051
		1052
		1053
	Yi Zhang, Tao Ge, and Xu Sun. 2020. Parallel data augmentation for formality style transfer . In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> , pages 3221–3228, Online. Association for Computational Linguistics.	1054
		1055
		1056
		1057
		1058
	Zheming Zhu, Delphine Bernhard, and Iryna Gurevych. 2010. A monolingual tree-based translation model for sentence simplification . In <i>Proceedings of the 23rd International Conference on Computational Linguistics (Coling 2010)</i> , pages 1353–1361, Beijing, China. Coling 2010 Organizing Committee.	1059
		1060
		1061
		1062
		1063
		1064
	A More Experimental Results	1065
	A.1 Out-of-domain Style Transfer	1066
	Table 5 shows automatic evaluation results of the ‘out-of-domain’ style transfer experiments, in	1067
		1068

Approach	CDS				GYFAC			
	TSS	MS	F	Agg.	TSS	MS	F	Agg.
GPT zero-shot	0.246 [‡]	0.657 [‡]	0.855 [‡]	0.138 [‡]	0.672 [‡]	0.752 [†]	0.909	0.455 [‡]
GPT 5-shot	0.289 [‡]	<u>0.708</u> [†]	0.868 [†]	0.175 [‡]	0.722 [‡]	<u>0.752</u> [†]	<u>0.902</u>	0.486 [‡]
STRAP	0.426 [‡]	0.629 [‡]	0.810 [‡]	0.194 [‡]	0.692 [‡]	0.689 [‡]	<u>0.852</u> [‡]	0.402 [‡]
STEER	<u>0.654</u> [†]	0.706 [†]	0.927	<u>0.426</u> [†]	<u>0.850</u> [†]	0.734 [‡]	0.875	<u>0.544</u> [†]
ASTRAPOP	0.579 [‡]	0.606 [‡]	0.808 [‡]	0.259 [‡]	0.816 [†]	0.685 [‡]	0.863 [‡]	<u>0.479</u> [‡]
STAMP	0.787	0.816	<u>0.877</u> [†]	0.562	0.964	0.864	0.827 [‡]	0.687

Table 5: The automatic evaluation results on out-of-domain inputs on the CDS and the GYAFAC datasets. The best and the 2nd best scores in each column are shown in **bold** and underline, respectively. “†” and “‡” indicate the score is significantly ($p < 0.05$) worse than the best score and the top 2 scores in the same column, respectively, determined by resampling t-test.

1069 which we transfer the texts in each dataset to the
1070 styles in the other dataset, in order to determine
1071 whether our results hold up when transferring be-
1072 tween styles of different provenance. They do; the
1073 out-of-domain results are generally consistent with
1074 the in-domain results. The best model in each col-
1075 umn in Table 5 is the same as Table 1, which is also
1076 true for the second best model in most columns.
1077 Also, STAMP still has the best TSS, MS, and
1078 aggregated score (Agg.) among all approaches,
1079 and STEER still has the best overall performance
1080 (Agg.) among the baselines.

1081 A.2 More Human Evaluation Results

Approach	TSS _h	MS _h	F _h
GPT 5-shot	0.59	<u>1.48</u>	<u>1.79</u>
STEER	0.69	1.24 [‡]	1.84
STAMP	<u>0.64</u>	1.48	1.57 [‡]

Table 6: Raw human evaluation scores on in-domain inputs on the CDS datasets. The best and 2nd best scores in each column are shown in **bold** and underline, respectively. “‡” indicates a statistically significant difference ($p < 0.05$) between the top two models determined by independent sample t-test. No significant difference is found in any other model pairs.

1082 The raw scores from the human evaluation and
1083 the result of the t-test are shown in Table 6. No
1084 significant difference is found between any model
1085 pairs in TSS_h¹⁷, but MS_h and F_h are generally
1086 consistent with the automatic evaluation results.
1087 Specifically, STAMP and GPT 5-shot transfer are
1088 significantly better than STEER in meaning simi-

¹⁷which is expected since style classification is difficult for human annotators (Krishna et al., 2020; Hallinan et al., 2023a).

1089 larity (MS), and STEER and GPT 5-shot transfer
1090 are significantly better than STAMP in fluency (F).

1091 A.3 Case Study

1092 We show an example from the CDS test set in Ta-
1093 ble 7 as a case study. In this example, we transfer
1094 a text in the style of music lyrics into the style
1095 of Shakespeare using STAMP and all baseline ap-
1096 proaches. STAMP maximally preserves the mean-
1097 ing of the original sentence and accurately reflects
1098 the target style using the words "tis" and "o'er".
1099 Other approaches either fail to generate strong tar-
1100 get style indicators or change the meaning of the
1101 original sentence to some extent.

1102 B More Implementation Details

1103 B.1 Statistical Significance Test

1104 We conduct a resampling paired t-test for the auto-
1105 matic evaluation results and an independent t-test
1106 for the human evaluation results. For the resam-
1107 pling paired t-test, we randomly select 10 subsets
1108 of 100 samples from the test set and perform a
1109 paired t-test on the mean scores of the subsets be-
1110 tween each pair of models. For the independent
1111 t-test, we use all available samples from the human
1112 study without resampling.

1113 B.2 Hyperparameters

1114 We sample same-sized training and validation sub-
1115 sets for CDS and GYAFAC, and use the same hyper-
1116 parameters to train STAMP on the two datasets to
1117 reduce the cost for more hyperparameter searching.
1118 We list all hyperparameters for STAMP in Table 8,
1119 Table 9, Table 10, Table 11, and Table 12.

Original	I'm not sorry that it's over
GPT zero-shot	I doth not lament the end that hath befallen us.
GPT 5-shot	I doth not lament the ending of this tale.
STRAP	I am not sorry That he is gone.
ASTRAPOP	Now is the winter of our discontent Made glorious summer by this sun of York.
STEER	I do not regret that it is done.
STAMP	I am not sorry That 'tis o'er.

Table 7: A style transfer example from the style of music lyrics to the style of Shakespeare.

Parameter	f_{cls}	f_{para}	$f_{p \rightarrow t}$	$f_{s \rightarrow t}$
learning rate	5e-5	5e-5	5e-5	5e-5
batch size	32	32	8	16
# epochs	6	10	6	12

Table 8: Training hyperparameters for all supervised fine-tuned models.

Parameter	f_{PO}
learning rate	2e-6
β	0.1
batch size	32
# epochs	16
k_{PO}	10
N_{iter}	10

Table 9: Training hyperparameters for iterative preference optimization.

Parameter	
target modules	q_proj, v_proj
rank	16
α	32
dropout	0.05

Table 10: LoRA Hyperparameters.

Parameter	$D_{p \rightarrow t}$	$D_{s \rightarrow t}$	D_{PO}
top p	1.0	1.0	1.0
temperature	0.5	0.7	1.0
$k_{para/sft/po}$	20	90	10
$\tau_{textMS/max}$	-	8	6

Table 11: Generation hyperparameters for dataset construction.

Parameter	Evaluation
top p	1.0
temperature	0.7

Table 12: Generation hyperparameters for dataset evaluation.

B.3 GPT prompt templates

1120

We elaborate on the prompts used for GPT zero- and 5-shot style transfer on CDS and GYAFC in Table 13 and Table 14, respectively.

1121

1122

1123

B.4 Hardware and Runtime

1124

We train all components of STAMP using Nvidia A40-48GB GPUs. The number of GPUs and time used to train each model on each dataset are shown in Table 15.

1125

1126

1127

1128

B.5 Human Evaluation Instructions

1129

The instructions used in the human evaluation for all three objectives are shown in Table 17 including the questions asked and the detailed explanation for each level in the Likert scale.

1130

1131

1132

1133

C Scientific Artifacts

1134

C.1 Use of Existing Artifacts

1135

The existing artifacts used in this work and their licenses are listed in Table 16. Our use of the existing artifacts is consistent with their intended use specified by their licenses.

1136

1137

1138

1139

C.2 Created Artifacts

1140

We create a new text style transfer training framework, STAMP, and release the code under the MIT license. Considering ethical implications, STAMP is only intended for research purposes, which is compatible with the original access conditions of all existing artifacts used in STAMP.

1141

1142

1143

1144

1145

1146

Zero-shot	Rewrite the following sentence into the style of [target style]. Original Sentence: [input text] Rewritten Sentence:
5-shot	Here are some examples of sentences in the style of [target style]: [example 1] [example 5] Rewrite the following sentence into the style of [target style]. Original Sentence: [input text] Rewritten Sentence:

Table 13: GPT zero- and 5-shot prompts for style transfer on CDS.

Zero-shot	Rewrite the following sentence in a(n) (in)formal style. Original Sentence: [input text] Rewritten Sentence:
5-shot	Here are some examples of sentences in a(n) (in)formal style: [example 1] [example 5] Rewrite the following sentence in a(n) (in)formal style. Original Sentence: [input text] Rewritten Sentence:

Table 14: GPT zero- and 5-shot prompts for style transfer on GYAFC.

	ParaNMT	CDS				GYAFC			
	f_{para}	f_{cls}	$f_{p \rightarrow t}$	$f_{s \rightarrow t}$	f_{PO}	f_{cls}	$f_{p \rightarrow t}$	$f_{s \rightarrow t}$	f_{PO}
# GPUs (A40s)	×2	×2	×2	×2	×4	×2	×2	×2	×2
Times (hrs)	3.4	0.4	1.1	1.0	35.2	0.1	0.2	0.2	7.4

Table 15: Training hardware and runtime for each component in STAMP on CDS and GYAFC.

Type	Name	License
Dataset	CDS: Corpus of Diverse Styles	MIT
	GYAFC: Grammarly’s Yahoo Answers Formality Corpus	Custom (research-only)
Model	LLaMA-2-7B (6.7B)	Meta
	GPT-3.5-turbo-0125 (-)	MIT
	RoBERTa-large (355M)	MIT
	RoBERTa-large CoLA Classifier (355M)	MIT
	SBERT all-mpnet-base-v2 (109M)	Apache-2.0
Library	Transformers	Apache-2.0
	PEFT	Apache-2.0
	TRL	Apache-2.0
	Sentence Transformers	Apache-2.0

Table 16: Datasets, models, and software libraries used in this work. The number of parameters of each model is indicated in the parentheses next to the model name.

TSS_h	Question	Based on the examples above, what is the style of the following text?
	Similar	Most of the meaning (75% or more) of the two passages is the same.
	Somewhat Similar	Large portions (50-75%) of the passages are the same, but there are significant sections that differ or are present in only one passage.
	Not Similar	Only small portions (less than 50%) of the passages are the same.
MS_h	Question	How similar are the following two texts?
	Fluent	Very clear, grammatical english (need not be formal); the meaning of the sentence is well understood. A small number of errors are ok.
	Somewhat Fluent	There are grammatical errors, possibly numerous, but the meaning can be understood.
	Not Fluent	The grammatical errors make it very difficult to understand the meaning.
F_h	Question	How fluent is the following text?

Table 17: Instructions used in the human evaluation.