

Balancing the Style-Content Trade-Off in Sentiment Transfer Using Polarity-Aware Denoising

Sourabrata Mukherjee^[0000-0002-1713-2769], Zdeněk Kasner^[0000-0002-5753-5538], and
Ondřej Dušek^[0000-0002-1415-1702]

Charles University, Faculty of Mathematics and Physics,
Institute of Formal and Applied Linguistics
Prague, Czechia
{mukherjee, kasner, odusek}@ufal.mff.cuni.cz

Abstract. Text sentiment transfer aims to flip the sentiment polarity of a sentence (positive to negative or vice versa) while preserving its sentiment-independent content. Although current models show good results at changing the sentiment, content preservation in transferred sentences is insufficient. In this paper, we present a sentiment transfer model based on polarity-aware denoising, which accurately controls the sentiment attributes in generated text, preserving the content to a great extent and helping to balance the style-content trade-off. Our proposed model is structured around two key stages in the sentiment transfer process: better representation learning using a shared encoder and sentiment-controlled generation using separate sentiment-specific decoders. Empirical results show that our methods outperforms state-of-the-art baselines in terms of content preservation while staying competitive in terms of style transfer accuracy and fluency. Source code, data, and all other related details are available on Github.¹

Keywords: Sentiment Transfer · Text Style Transfer · Natural Language Generation

1 Introduction

Text sentiment transfer is the task of changing the sentiment polarity of a text while retaining sentiment-independent semantic content (e.g., “The food was tasteless” to “The food was delicious”) [26, 20, 12, 14]. It has been introduced in the context of textual style transfer, where positive and negative sentiment are considered distinct styles [14]. Style transfer is motivated by various writing assist tasks for copywriting or personalized chatbots, e.g. changing review sentiment, debiasing or simplifying a news text, or removing offensive language [14, 25, 9].

With the success of deep learning in the last decade, a variety of neural methods have been proposed for this task [27, 9]. If parallel data are provided, standard sequence-to-sequence models can be directly applied [23]. However, due to lack of parallel corpora (paired sentences with opposite sentiment and otherwise identical content), learning sentiment transfer – and text style transfer in general – from unpaired data represents a substantial research challenge.

¹ <https://github.com/SOURO/polarity-denoising-sentiment-transfer>

A first approach to learning text style transfer from unpaired data disentangles text representation in a latent space into style-independent content and stylistic attributes (such as sentiment polarity) and applies generative modeling [7, 26, 20]. The latent representation needs to preserve the meaning of the text while abstracting away from its stylistic properties, which is not trivial [11]. In fact, disentanglement is impossible in theory without inductive biases or other forms of supervision [13]. A second line of research is prototype editing [12, 3], which focuses specifically on style marker words (also called *pivot* words, e.g. sentiment polarity indicating words such as “good” or “bad”). The approach typically extracts a sentence “template” without the pivots and then fills in pivots marking the target style. However, since the pivot words are typically extracted using simple unsupervised probabilistic methods, they are difficult to distinguish from content words, which again leads to content preservation errors.

Our work combines both research branches and extends them, using additional supervision. The supervision comes from a sentiment dictionary, which is applied on pivot words within the context of generative models to learn better latent representations via the task of polarity-aware denoising. First, we randomly delete (or mask) pivot word(s) of input sentences. Then a shared encoder pre-trained on general domain helps in preparing a latent representation, followed by separate sentiment-specific decoders that are used to change the sentiment of the original sentence. We follow back-translation for style transfer approach proposed by Prabhunoye et al. [20] to represent the sentence meaning in the latent space.

Our contributions are summarized as follows:

- We design a sentiment transfer model using an extended transformer architecture and polarity-aware denoising. Our use of separate sentiment-specific decoders and polarity-aware denoising training allow more control over both the target sentiment and the sentiment-independent content.
- We derive a new non-parallel sentiment transfer dataset from the Amazon Review Dataset [17]. It is more topically diverse than earlier used datasets Yelp [12] and IMDb [2], which were majorly focused on movie and restaurant/business-related reviews. Our dataset and code is publicly available.¹
- We introduce polarity-masked BLEU (MaskBLEU) and similarity score (MaskSim) for automatic evaluation of content preservation in this task. These metrics are derived from the traditional BLEU score [19] and Sentence BERT-based cosine similarity score [24]. In our approach, we mask polarity words beforehand for sentiment-independent content evaluation.
- Both automatic and human evaluations on our dataset show that our proposed approach generally outperforms state-of-the-art (SotA) baselines. Specifically, with respect to content preservation, our approach achieves substantially better performance than other methods.

2 Related Work

Sentiment Transfer A common approach to the sentiment transfer task is to separate content and style in a latent space, and then adjust the separated style. Hu et al. [7] use a variational autoencoder and factor its latent representation into a style-independent and

stylistic parts. Fu et al. [4] compare a multi-decoder model with a setup using a single decoder and style embeddings. Shen et al. [26] propose a cross-aligned auto-encoder with adversarial training. Prabhume et al. [20] propose to perform text style transfer through back-translation. In a recent work, He et al. [6] apply variational inference. Although these approaches successfully change the text style, they also change the text content, which is a major problem. Many previous methods [7, 26, 4, 20] formulate the style transfer using the encoder-decoder framework. The encoder maps the text into a style-independent latent representation, and the decoder generates the target text using the latent representation and a style marker. Again, a major issue of these models is poor preservation of non-stylistic semantic content.

Content Preservation To further deal with the above problem, Li et al. [12] first extract content words by deleting phrases, then retrieve new phrases associated with the target attribute, and finally use a neural model to combine these into a final output. Luo et al. [14] employ a dual reinforcement learning framework with two sequence-to-sequence models in two directions, using style classifier and back-transfer reconstruction probability as rewards. Though these works show some improvement, they are still not able to properly balance preserving the content with transferring the style. Our polarity-aware denoising technique aims to solve this problem by specifically targeting and changing polarity words while preserving the rest of the content (see Section 4.3).

Evaluation Another challenge remains in the evaluation of textual style transfer. Previous work on style transfer [7, 20, 2, 6] has re-purposed evaluation metrics from other fields, such as BLEU from machine translation [19] and PINC from paraphrasing [1]. However, these metric cannot evaluate style transfer specifically with respect to preservation of content [27] as they do not take into account the necessity of changing individual words when altering the style. Intended differences between the source sentence and the transferred sentence are thus penalized. In this regard, we have introduced polarity masked BLEU score (MaskBLEU) and polarity masked similarity measure (MaskSim), where we have masked the polarity words beforehand (see Section 5.3).

3 Approach

3.1 Task Definition

Given two datasets, $X_{pos} = \{x_1^{(pos)}, \dots, x_m^{(pos)}\}$ and $X_{neg} = \{x_1^{(neg)}, \dots, x_n^{(neg)}\}$ which represent two different sentiments *pos* and *neg*, respectively, our task is to generate sentences of the desired sentiment while preserving the meaning of the input sentence. Specifically, we alter samples of dataset X_{pos} such that they belong to sentiment *neg* and samples of X_{neg} such that they belong to sentiment *pos*, while sentiment-independent content is preserved. We denote the output of dataset X_{pos} transferred to sentiment *neg* as $X_{pos \rightarrow neg} = \{\hat{x}_1^{(neg)}, \dots, \hat{x}_m^{(neg)}\}$ and the output of dataset X_{neg} transferred to sentiment *pos* as $X_{neg \rightarrow pos} = \{\hat{x}_1^{(pos)}, \dots, \hat{x}_n^{(pos)}\}$.

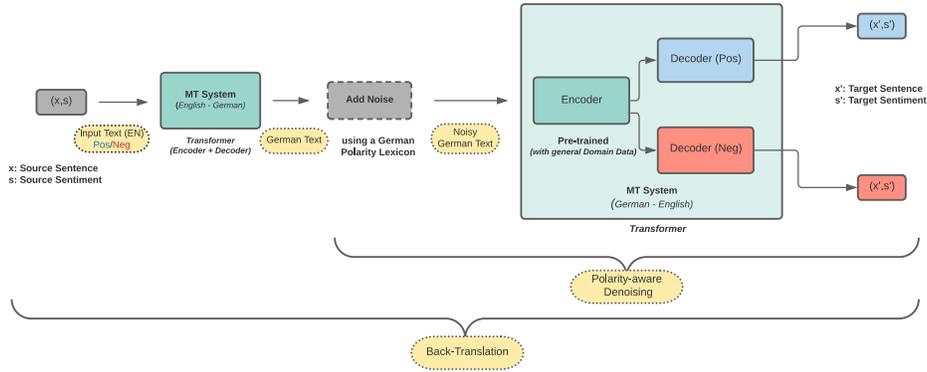


Fig. 1. Our sentiment transfer pipeline. In the pipeline, we (1) *translate* the source sentence from English to German using a transformer-based machine translation (MT) system; (2) *apply noise* on the German sentence using a German polarity lexicon; (3) *encode* the German sentence to latent representation using an encoder of German-to-English translation model; (4) *decode* the shared latent representation using the decoder for the opposite sentiment.

3.2 Model Overview

Figure 1 shows the overview of our proposed architecture. Following Prabhumoye et al. [20], we first translate the input text x in the base language to a chosen intermediate language \bar{x} using a translation model (see Section 4.1).² Next, we prepare a noisy text x_{noise} from \bar{x} using polarity-aware noising with word deletion or masking probabilities θ_N (see Section 4.3):

$$x_{noise} = Noise(\bar{x}; \theta_N). \quad (1)$$

We provide x_{noise} to the encoder of the $\bar{x} \rightarrow \hat{x}$ back-translation model (where \hat{x} is a text in the base language with changed sentiment polarity). The encoder first converts the text to the latent representation z as follows:

$$z = Encoder(x_{noise}; \theta_E), \quad (2)$$

where θ_E represent the parameters of the encoder.

Two separate sentiment-specific decoders are trained to decode the original positive and negative inputs by passing in their latent representations z :

$$x_{pos} = Decoder_{pos}(z; \theta_{D_{pos}}) \quad (3)$$

$$x_{neg} = Decoder_{neg}(z; \theta_{D_{neg}}). \quad (4)$$

At inference time, sentiment transfer is achieved by decoding the shared latent representation using the decoder trained for the opposite sentiment, as follows:

$$\hat{x}_{neg} = Decoder_{pos}(z; \theta_{D_{pos}}) \quad (5)$$

$$\hat{x}_{pos} = Decoder_{neg}(z; \theta_{D_{neg}}) \quad (6)$$

² We work with English as base language and German as intermediate language, see Section 5.1.

where \hat{x}_{neg} , \hat{x}_{pos} are the sentences with transferred sentiment conditioned on z and $\theta_{D_{pos}}$ and $\theta_{D_{neg}}$ represent the parameters of the positive and negative decoders, respectively.

4 Model Variants

In all our experiments, we train sentiment transfer models using back-translation (Section 4.1) based on the transformer architecture [28]. First, we present baselines for sentiment transfer with simple style conditioning (Section 4.2). Next, we propose an approach based on an extended transformer architecture where we use separate modules (either the whole transformer model, or the transformer decoder only) for the respective target sentiment (Section 4.2). We further improve upon our approach using polarity-aware denoising (Section 4.3), which we propose as a new scheme for pre-training the sentiment transfer models.

4.1 Back-translation

Back-translation for style transfer was introduced in Prabhumoye et al. [20]. Following their approach, we use translation into German and subsequent encoding in a back-translation model to get a latent text representation for our sentiment transfer task. Prior work has also shown that the process of translating a sentence from a source language to a target language retains the meaning of the sentence but does not preserve the stylistic features related to the author’s traits [21]. A pure back-translation approach (without any specific provisions for sentiment) is referred to as *Back-Translation* in our experiments.

We also experimented with an auto-encoder, but we have found that the back-translation model gives better results for sentiment transfer. We hypothesise that it is due to the fact that back-translation prevents the system from sticking to a particular wording, resulting in a more abstract latent representation.

4.2 Our Baseline Models

In addition to a pure back-translation model, we present several straightforward baselines:

- *Style Tok* is a back-translation model with added sentiment identifiers ($\langle pos \rangle$ or $\langle neg \rangle$) as output starting tokens. At the time of sentiment transfer, we decode the output with a changed sentiment identifier ($\langle pos \rangle \rightarrow \langle neg \rangle$, $\langle neg \rangle \rightarrow \langle pos \rangle$).
- *Two Sep. transformers*: To get more control over sentiment-specific generation, we train two separate transformer models for positive and negative sentiment, using only sentences of the respective target sentiment. During inference, the model is fed with inputs of the opposite sentiment, which it did not see during training.
- *Shrd Enc + Two Sep Decoders*: We extend the above approach by keeping decoders separate, but using a shared encoder. During training, all examples are passed through the shared encoder, but each decoder is trained to only generate samples of one sentiment. Sentiment transfer is achieved by using the decoder for the opposite sentiment.
- *Pre Training Enc*: Following Gururangan et al. [5], we introduce a variant where the shared encoder is pretrained for back-translation on general-domain data. The pre-trained encoder is then further fine-tuned during sentiment transfer training.

4.3 Polarity-Aware Denoising

We devise a task-specific pre-training [5] scheme for improving the sentiment transfer abilities of the model. The idea of our pre-training scheme—*polarity-aware denoising*—is to first introduce noise, i.e. delete or mask a certain proportion of words in the intermediate German input to the back-translation step, then train the model to remove this noise, i.e. produce the original English sentence with no words deleted or masked. To decide which words get deleted or masked, we use automatically obtained sentiment polarity labels (see Section 5.2 for implementation details). This effectively adds more supervision to the task on the word level. We apply three different approaches: deleting or masking (1) *general* words (i.e., all the words uniformly), (2) *polarity* words (i.e., only high-polarity words according to a lexicon), or (3) both *general* and *polarity* words (each with a different probability).

We use polarity-aware denoising during encoder pretraining, following the shared encoder and separate decoders setup from Section 4.2. The encoder is further fine-tuned during the sentiment transfer training.

5 Experiments

We evaluated and compared our approaches described in Section 4 with several state-of-the-art systems [26, 20, 12, 14, 29, 6] on two datasets (see Section 5.1). We first describe our training setup (Section 5.2), then detail our automatic evaluation metrics (Section 5.3) and human evaluation (Section 5.4), and finally discuss the results (Section 5.5).

5.1 Datasets

For our back-translation process and model pretraining, we used the WMT14 English-German (*en-de*) dataset (1M sentences) from Neidert et al. [16].

For finetuning and experimental evaluation, we built a new English dataset for sentiment transfer, based on the Amazon Review Dataset [17]. We have selected Amazon Review because it is more diverse topic-wise (books, electronics, movies, fashion, etc.) than existing sentiment transfer datasets, Yelp [12] and IMDb [2], which are majorly focused on movie and restaurant/business-related reviews. For comparison with previous work, we also evaluate our models on the benchmark IMDb dataset [2].

While the Amazon Review data is originally intended for recommendation, it lends itself easily to our task. We have split the reviews into sentences and then used a pre-trained transformer-based sentiment classifier [30] to select sentences with high polarity. Our intuition is that high-polarity sentences are more informative for the sentiment transfer task than neutral sentences. We filter out short sentences (less than 5 words) since it is hard to evaluate content preservation for these sentences. We also ignored sentences with repetitive words (e.g., "*no no no no thanks thanks.*") because these sentences are noisy and do not serve as good examples for the sentiment transfer model. We aim at comparable size to existing datasets [12]: the resulting data has 102k positive and 102k negative examples in total, with 1+1k reserved for validation and testing for each sentiment. The average sentence length in our data is 13.04 words.

5.2 Training Setup

In all our experiments, we use a 4-layer transformer [28] with 8 attention heads per layer. Both embedding and hidden layer size are set to 512. The same model shape is used for both the initial translation into German and the subsequent model handling back-translation, denoising, and sentiment transfer.

We use a German polarity lexicon to automatically identify pivot words for polarity-aware denoising. We prepared the German polarity lexicon by first translating the words from German to English using an off-the-shelf translation system [10], followed by labeling the words with *positive* and *negative* labels using the English NLTK Vader lexicon [8]. We performed a manual check of the results on a small sample.

We test various combinations of noise settings w.r.t. noise probability, noise type (general or polarity-aware denoising), and noise mode (deleting or masking). Parameter values are pre-selected based on our preliminary experiments with the translation model (see Section 4.1). The parameters are encoded in the name of the model as used in Table 1 (see the table caption for details).

5.3 Automatic Evaluation Metrics

To evaluate the performance of the models, we compare the generated samples along three different dimensions using automatic metrics, following previous work: (1) style control, (2) content preservation, and (3) fluency.

Standard Metrics

- *Style Accuracy*: Following prior work, we measure sentiment accuracy automatically by evaluating the sentiment of transferred sentences. We use a pre-trained transformer-based sentiment analysis pipeline³ from Huggingface [30].
- *Fluency*: We use the negative log-likelihood from the GPT-2 [22] language model as an indirect metric for evaluating fluency. For context, we also calculate average sentence lengths of the sentiment-transferred sentences.
- *Content Preservation*: Following previous work, we compute BLEU score [19] between the transferred sentence and its source. Higher BLEU indicates higher n-gram overlap between the sentences, which is generally viewed as proxy for content preservation. We also compute Sentence BERT [24] based cosine similarity score to match the vector space semantic similarity between the source and the transferred sentence. None of the techniques is capable of evaluating style transfer methods specifically with respect to preservation of content in style transfer [27]. These metrics do not take into account the necessity of changing individual words while altering the sentence style. Intended differences between the source sentence and the transferred sentence are thus penalized.

Newly Introduced Metrics for Content Preservation To avoid the problems of the commonly used metrics, it makes sense in sentiment transfer to evaluate the content

³ <https://huggingface.co/distilbert-base-uncased-finetuned-sst-2-english>

Table 1. Automatic evaluation. *Accuracy*: Sentiment transfer accuracy. *Sim* and *B*: Cosine similarity and BLEU score between input and sentiment-transferred sentence. *M/Sim* and *M/B*: MaskSim and MaskBLEU (similarity and BLEU with polarity words masked, see Section 5.3). *LM*: Average log probability assigned by vanilla GPT-2 language model. *Avg*: Average length of transferred sentences. *Avg*: Average of sentiment transfer accuracy, 100*MaskSim and MaskBLEU. Scores are based on a single run, with identical random seeds. First two sections show our own baselines, third section shows our models with denoising (with the best settings denoted SCT₁ and SCT₂, see Section 5.5). The bottom section shows a comparison with state-of-the-art models. Names of models with denoising reflect settings as follows: *W* denotes WMT pretraining data, *A* denotes Amazon finetuning data; the following tokens denote noise probability values are associated with the respective data. *G/P* represents general/polarity token noising, *D/M* represents noising mode deletion/masking. E.g. *WG03P08-AG03P08-D*: noise probabilities on WMT and Amazon data are identical, noising by deletion on both general and polarity token noising is applied (with probabilities 0.3 and 0.8, respectively).

Models	Acc	Sim	M/Sim	B	M/B	LM	Len	Avg
Back-Translation Only (Section 4.1)								
<i>Back-translation only</i>	0.4	0.828	0.768	28.0	45.3	-78.6	11.9	40.9
Our Baseline Models (Section 4.2)								
<i>Style Tok</i>	13.2	0.536	0.560	4.8	8.6	-52.1	7.6	25.9
<i>Two Sep. transformers</i>	89.3	0.394	0.611	6.8	19.6	-79.0	13.7	56.7
<i>Shrd Enc + Two Sep Decoders</i>	88.1	0.397	0.600	7.3	20.1	-78.0	12.5	56.0
<i>Pre Training Enc</i>	55.3	0.592	0.732	22.6	33.9	-93.3	13.4	54.1
Our Models (with Denoising) (Section 4.3)								
<i>WG01-AG01-D</i>	71.4	0.517	0.694	17.1	29.8	-88.7	13.7	56.9
<i>WG01-AG01-M</i>	68.0	0.536	0.711	19.4	31.1	-86.3	12.6	56.7
<i>WG03-AG03-D</i>	83.0	0.447	0.648	11.7	24.4	-83.0	13.7	57.4
<i>WG03-AG03-M</i>	78.8	0.481	0.669	14.2	28.2	-82.7	13.0	58.0
<i>WP08-AP08-D</i>	66.9	0.528	0.701	19.5	31.3	-82.8	12.4	56.1
<i>WP08-AP08-M</i>	64.0	0.547	0.726	21.4	34.0	-89.1	12.9	56.9
<i>WPI-API-D</i>	58.7	0.570	0.727	22.7	33.1	-87.2	12.2	54.8
<i>WPI-API-M</i>	58.9	0.567	0.716	22.2	33.0	-86.5	12.2	54.5
<i>WG03-AG01-D</i>	68.0	0.529	0.697	17.9	30.9	-89.5	13.3	56.2
<i>WG03-AG01-M</i>	80.7	0.473	0.665	13.9	27.5	-82.8	13.1	58.2
<i>WG01-AG03-D (=SCT₂)</i>	85.2	0.441	0.646	11.8	25.4	-79.8	13.1	58.4
<i>WG01-AG03-M</i>	70.0	0.534	0.711	19.7	32.3	-84.3	12.4	57.8
<i>WP08-API-D</i>	61.6	0.578	0.736	22.5	35.0	-94.4	13.4	56.7
<i>WP08-API-M</i>	60.9	0.554	0.724	22.0	33.3	-85.5	12.6	55.6
<i>WPI-AP08-D</i>	68.5	0.525	0.699	19.3	31.1	-84.0	12.4	56.5
<i>WPI-AP08-M</i>	61.1	0.560	0.714	21.5	32.9	-86.0	12.1	55.1
<i>WG03-AP08-D</i>	67.0	0.533	0.697	20.3	31.7	-84.3	12.5	56.1
<i>WG03-AP08-M</i>	65.7	0.546	0.725	21.2	33.5	-85.0	12.5	57.2
<i>WP08-AG03-D</i>	83.3	0.436	0.635	11.0	24.3	-80.5	13.3	57.0
<i>WP08-AG03-M</i>	79.6	0.473	0.665	13.2	26.9	-83.1	13.2	57.6
<i>WG03P08-AG03P08-D</i>	65.5	0.547	0.705	20.3	32.6	-90.4	13.2	56.2
<i>WG03P08-AG03P08-M (=SCT₁)</i>	82.0	0.460	0.665	13.7	27.4	-79.6	12.8	58.6
State-of-the-Art Models								
Shen et al. [26]	88.6	0.346	0.513	3.2	18.3	-74.0	10.9	52.7
Li et al. [12]	69.9	0.457	0.632	14.7	25.3	-85.1	12.2	52.8
Luo et al. [14]	92.4	0.279	0.468	0.0	9.1	-42.0	7.8	49.4
Prabhumoye et al. [20]	93.5	0.308	0.504	0.9	15.2	-61.0	10.3	53.0
Wang et al. [29]	79.3	0.385	0.545	10.6	20.3	-116.8	15.1	51.4
He et al. [6]	91.5	0.352	0.542	9.5	21.8	-65.9	8.2	55.8

Table 2. Automatic evaluation on the IMDb Dataset (see Table 1 for metrics explanation).

Models	Acc	Sim	M/Sim	B	M/B	LM	Len	Avg
Prabhumoye et al. [20]	87.1	0.345	0.480	2.7	14.3	-63.5	10.0	49.8
Li et al. [12]	21.0	0.587	0.668	18.3	25.9	-83.6	15.3	37.9
Wang et al. [29]	84.0	0.357	0.456	9.2	13.2	-63.9	10.8	47.6
He et al. [6]	81.7	0.458	0.576	29.0	41.8	-83.6	15.3	60.4
SCT ₁ (WG03P08-AG03P08-M)	85.3	0.435	0.612	28.6	42.3	-86.4	15.9	62.9
SCT ₂ (WG01-AG03-D)	88.2	0.379	0.588	25.8	39.2	-79.6	15.1	62.1

Table 3. Human evaluation of sentiment transfer quality, content preservation, and fluency. Average of 1-5 Likert scale ratings on 100 examples from our Amazon Review data.

Models	Sentiment	Content	Fluency
Prabhumoye et al. [20]	3.95	1.19	3.56
Li et al. [12]	3.35	2.3	3.34
Wang et al. [29]	3.48	1.67	2.54
He et al. [6]	3.69	1.66	3.26
SCT ₁ (WG03P08-AG03P08-M)	3.94	2.61	3.73
SCT ₂ (WG01-AG03-D)	3.99	2.56	3.79

and similarity while ignoring any polarity tokens. Thus, we introduce MaskBLEU and MaskSim scoring methods – these are identical to BLEU and cosine similarity, but they are computed on sentences where pivot words (based on NLTK Vader sentiment dictionary [8]) have been masked. This allows measuring content preservation while ignoring the parts of the sentences that need to be changed.

5.4 Human Evaluation

As automated metrics for language generation do not correlate well with human judgements [18], we conduct an in-house human evaluation with five expert annotators. We randomly select 100 sentences (50 for each sentiment) from the our Amazon Review test set. The annotators rate model outputs on using a 1-5 Likert scale for style control, content preservation and fluency.

5.5 Results

Automatic Metrics results on our Amazon Review data are shown in Table 1. Overall, there is clearly a tradeoff between preserving sentiment-independent content and achieving the desired target sentiment. Models which perform very well in sentiment transfer usually achieve worse results on content preservation. This tradeoff is documented by correlations between the automatic metrics (see Figure 2). Sentiment accuracy is negatively correlated with BLEU score, similarity measures as well as our newly introduced MaskBLEU and MaskSim scores.

The translation-only and style token baselines do not perform well on changing the sentiment. Using two separate decoders leads to major sentiment transfer improvements,

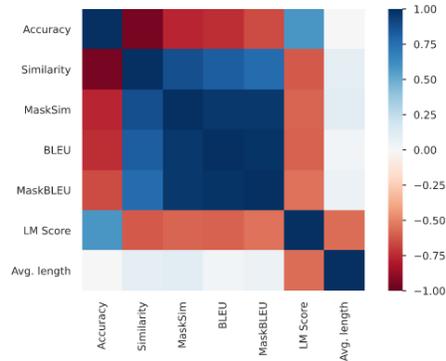


Fig. 2. Correlations between automatic evaluation metrics on our Amazon Review data: sentiment accuracy is negatively correlated with BLEU, semantic similarity, and their masked variants.

but content preservation is poor. Using the pre-trained encoder has helped to improve the content preservation, but sentiment transfer accuracy degrades significantly.

The main motivation for our work was to find a denoising strategy which offers the best balance between sentiment transfer and content preservation. Our results suggest putting an emphasis on denoising high-polarity words results in the best ratio between the sentiment transfer accuracy and content preservation metrics. Additionally, our models show the ability to produce fluent sentences, as shown by the language model score: our models’ scores are similar to the back-translation baseline; other models only reach higher language model scores when producing very short outputs.

Overall, our denoising approaches are able to balance well between sentiment transfer and content preservation. On content preservation, they perform much better than state-of-the-art models, and they stay competitive on style accuracy. We selected two of our model variants – $SCT_1=WG03P08-AG03P08-M$ and $SCT_2=WG01-AG03-D$ – as the ones giving the best style-content trade-off (SCT) according to the average of sentiment accuracy, masked similarity and MaskBLEU (see Table 1).

Automatic metrics on the IMDb dataset [2] are shown in Table 2, comparing our selected SCT_1 and SCT_2 models with state-of-the-art. Our models outperform the state-of-the-art in terms of sentiment accuracy and reach competitive results in terms of similarity, BLEU, and fluency. Same as on our Amazon Review data, they provide the best style-content trade-off (according to the averaged metric defined in Table 1).

Human Evaluation Results: We compare our best SCT_1 and SCT_2 models (selected above) with four state-of-the-art models: two of the most recent models [29, 6], and the models with best accuracy [20] and MaskBLEU score [12].

We have evaluated over 600 model outputs. Results are presented in Table 3. The human evaluation results mostly agree with our automatic evaluation results. The results also show that our models are better in content preservation than the competitor models.

Table 4. Example outputs comparison on samples from our Amazon Reviews dataset. Sentiment marker words (pivots) are colored. Note that our models preserve content better than most others.

	Negative → Positive	Positive → Negative
Source	movie was a waste of money : this movie totally sucks .	my daughter loves them :)
Prabhumoye et al. [20]	stan is always a great place to get the food .	do n't be going here .
Li et al. [12]	our favorite thing was a movie story : the dream class roll !	my daughter said i was still not acknowledged .
Wang et al. [29]	movie is a delicious atmosphere of : this movie totally sucks movie !	i should not send dress after me more than she would said not ?
He et al. [6]	this theater was a great place , we movie totally amazing .	yup daughter has left ourselves .
SCT ₁ (WG03P08-AG03P08-M)	movie : a great deal of money : this movie is absolutely perfect .	my daughter hates it : my daughter .
SCT ₂ (WG01-AG03-D)	this movie is a great deal of money.	my daughter hated it .
Source	nothing truly interesting happens in this book .	best fit for my baby : this product is wonderful !!
Prabhumoye et al. [20]	very good for the best .	bad customer service to say the food , and it is n't .
Li et al. [12]	nothing truly interesting happens in this book .	my mom was annoyed with my health service is no notice .
Wang et al. [29]	nothing truly interesting happens in this book make it casual and spot .	do not buy my phone : this bad crap was worst than it ?
He et al. [6]	haha truly interesting happens in this book .	uninspired .
SCT ₁ (WG03P08-AG03P08-M)	in this book is truly a really great book .	not good for my baby : this product is great ! ! ! ! ! ! ! !
SCT ₂ (WG01-AG03-D)	in this book is truly awesome .	not happy for my baby : this product is not great ! !

We further examined a sample of the outputs in more detail to understand the behavior of different models. We found that state-of-the-art models tend to lose the content of the source sentence, as shown in the example outputs in Table 4. On the other hand, our models mostly preserve sentiment-independent content well while successfully transferring the sentiment. We conclude that with our models, there is a good balance between preserving the original sentiment-independent content and dropping the source sentiment, and existing state-of-the-art models tend to sacrifice one or the other.

6 Conclusions and Future Work

We proposed an approach for text sentiment transfer based on the transformer architecture and polarity-aware denoising. Experimental results on two datasets showed that our method achieves competitive or better performance compared to state-of-the-art. Our architecture provides a good style-content tradeoff mainly due to two elements:

(1) separate sentiment-specific decoders providing explicit target sentiment control, and (2) polarity-aware enhanced denoising removing sentiment implicitly at the token level. As shown by human evaluation and our manual inspection, our models still sometimes fail to preserve the meaning of the original. While we improve upon previous works on content preservation, this remains a limitation.

In the future, we plan to adapt our method to the different kind of style transfer tasks such as formality transfer or persona-based text generation. Lexicons for the required attribute makers can be extracted by mining stylistic markers from generic dictionaries, or from personality-annotated data [15]. We also intend to focus on better controlling content preservation with the use of semantic parsing.

Acknowledgments

This research was supported by Charles University projects GAUK 392221, GAUK 140320, SVV 260575 and PRIMUS/19/SCI/10, and by the European Research Council (Grant agreement No. 101039303 NG-NLG). It used resources provided by the LINDAT/CLARIAH-CZ Research Infrastructure (Czech Ministry of Education, Youth and Sports project No. LM2018101).

References

- [1] Chen, D., Dolan, W.B.: Collecting highly parallel data for paraphrase evaluation. In: Proc. ACL-HLT, pp. 190–200 (2011)
- [2] Dai, N., Liang, J., Qiu, X., Huang, X.: Style transformer: Unpaired text style transfer without disentangled latent representation. In: Proc. ACL, pp. 5997–6007 (2019)
- [3] Fu, Y., Zhou, H., Chen, J., Li, L.: Rethinking Text Attribute Transfer: A Lexical Analysis. In: Proc. INLG (2019)
- [4] Fu, Z., Tan, X., Peng, N., Zhao, D., Yan, R.: Style Transfer in Text: Exploration and Evaluation. In: Proc. AAAI, pp. 663–670 (2018)
- [5] Gururangan, S., Marasovic, A., Swayamdipta, S., Lo, K., Beltagy, I., Downey, D., Smith, N.A.: Don’t Stop Pretraining: Adapt Language Models to Domains and Tasks. In: Proc. ACL, pp. 8342–8360 (2020)
- [6] He, J., Wang, X., Neubig, G., Berg-Kirkpatrick, T.: A Probabilistic Formulation of Unsupervised Text Style Transfer. In: Proc. ICLR (2020)
- [7] Hu, Z., Yang, Z., Liang, X., Salakhutdinov, R., Xing, E.P.: Toward Controlled Generation of Text. In: Proc. ICML, vol. 70, pp. 1587–1596 (2017)
- [8] Hutto, C.J., Gilbert, E.: VADER: A Parsimonious Rule-Based Model for Sentiment Analysis of Social Media Text. In: Proc. ICWSM (2014)
- [9] Jin, D., Jin, Z., Hu, Z., Vechtomova, O., Mihalcea, R.: Deep Learning for Text Style Transfer: A Survey. Computational Linguistics pp. 1–51 (Dec 2021)
- [10] Kořarko, O., Variš, D., Popel, M.: LINDAT translation service (2019), URL <http://hdl.handle.net/11234/1-2922>
- [11] Lample, G., Subramanian, S., Smith, E.M., Denoyer, L., Ranzato, M., Boureau, Y.: Multiple-Attribute Text Rewriting. In: Proc. ICLR (2019)

- [12] Li, J., Jia, R., He, H., Liang, P.: Delete, Retrieve, Generate: a Simple Approach to Sentiment and Style Transfer. In: Proc. NAACL-HLT, pp. 1865–1874 (2018)
- [13] Locatello, F., Bauer, S., Lucic, M., Rätsch, G., Gelly, S., Schölkopf, B., Bachem, O.: Challenging Common Assumptions in the Unsupervised Learning of Disentangled Representations. In: Proc. ICML, vol. 97, pp. 4114–4124 (2019)
- [14] Luo, F., Li, P., Yang, P., Zhou, J., Tan, Y., Chang, B., Sui, Z., Sun, X.: Towards Fine-grained Text Sentiment Transfer. In: Proc. ACL, pp. 2013–2022 (2019)
- [15] Mairesse, F., Walker, M.A.: Controlling user perceptions of linguistic style: Trainable generation of personality traits. *Computational Linguistics* **37**(3), 455–488 (2011)
- [16] Neidert, J., Schuster, S., Green, S., Heafield, K., Manning, C.D.: Stanford University’s Submissions to the WMT 2014 Translation Task. In: Proc. WMT, pp. 150–156 (2014)
- [17] Ni, J., Li, J., McAuley, J.J.: Justifying Recommendations using Distantly-Labeled Reviews and Fine-Grained Aspects. In: Proc. EMNLP-IJCNLP, pp. 188–197 (2019)
- [18] Novikova, J., Dusek, O., Curry, A.C., Rieser, V.: Why We Need New Evaluation Metrics for NLG. In: Proc. EMNLP, pp. 2241–2252 (2017)
- [19] Papineni, K., Roukos, S., Ward, T., Zhu, W.: Bleu: a Method for Automatic Evaluation of Machine Translation. In: Proc. ACL, pp. 311–318 (2002)
- [20] Prabhume, S., Tsvetkov, Y., Salakhutdinov, R., Black, A.W.: Style Transfer Through Back-Translation. In: Proc. ACL, pp. 866–876 (2018)
- [21] Rabinovich, E., Patel, R.N., Mirkin, S., Specia, L., Wintner, S.: Personalized machine translation: Preserving original author traits. In: Proc. EACL, pp. 1074–1084 (2017)
- [22] Radford, A., Wu, J., Child, R., Luan, D., Amodei, D., Sutskever, I.: Language models are unsupervised multitask learners. Tech. report, Open AI (2019)
- [23] Rao, S., Tetreault, J.R.: Dear Sir or Madam, May I Introduce the GYAFC Dataset: Corpus, Benchmarks and Metrics for Formality Style Transfer. In: Proc. NAACL-HLT, pp. 129–140 (2018)
- [24] Reimers, N., Gurevych, I.: Sentence-BERT: Sentence Embeddings using Siamese BERT-Networks. In: Proc. EMNLP-IJCNLP, pp. 3980–3990 (2019)
- [25] Santos, C.N.d., Melnyk, I., Padhi, I.: Fighting offensive language on social media with unsupervised text style transfer. arXiv preprint arXiv:1805.07685 (2018)
- [26] Shen, T., Lei, T., Barzilay, R., Jaakkola, T.S.: Style Transfer from Non-Parallel Text by Cross-Alignment. In: Proc. NeurIPS, pp. 6830–6841 (2017)
- [27] Toshevska, M., Gievska, S.: A Review of Text Style Transfer using Deep Learning. *IEEE Transactions on Artificial Intelligence* (2021)
- [28] Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A.N., Kaiser, L., Polosukhin, I.: Attention is all you need. In: Proc. NeurIPS, pp. 5998–6008 (2017)
- [29] Wang, K., Hua, H., Wan, X.: Controllable unsupervised text attribute transfer via editing entangled latent representation. In: Proc. NeurIPS, pp. 11034–11044 (2019)
- [30] Wolf, T., Debut, L., Sanh, V., Chaumond, J., Delangue, C., Moi, A., Cistac, P., Rault, T., Louf, R., Funtowicz, M., Davison, J., Shleifer, S., von Platen, P., Ma, C., Jernite, Y., Plu, J., Xu, C., Scao, T.L., Gugger, S., Drame, M., Lhoest, Q., Rush, A.M.: Transformers: State-of-the-art natural language processing. In: Proc. EMNLP, pp. 38–45 (2020)