

Generative Sentiment Transfer via Adaptive Masking

Yingze Xie¹, Jie Xu¹, LiQiang Qiao¹, Yun Liu², Feiren Huang³, Chaozhuo Li⁴

¹ School of Information Science and Technology, Beijing Foreign Studies University, Beijing, China

² Moutai Institute, China

³ School of Information Science and Technology, Jinan University, China

⁴ Microsoft Research Asia, Beijing, China

Abstract. Sentiment transfer aims at revising the input text to satisfy a given sentiment polarity while retaining the original semantic content. The nucleus of sentiment transfer lies in precisely separating the sentiment information from the content information. Existing explicit approaches generally identify and mask sentiment tokens simply based on prior linguistic knowledge and manually-defined rules, leading to low generality and undesirable transfer performance. In this paper, we view the positions to be masked as the learnable parameters, and further propose a novel AM-ST model to learn adaptive task-relevant masks based on the attention mechanism. Moreover, a sentiment-aware masked language model is further proposed to fill in the blanks in the masked positions by incorporating both context and sentiment polarity to capture the multi-grained semantics comprehensively. AM-ST is thoroughly evaluated on two popular datasets, and the experimental results demonstrate the superiority of our proposal.

1 Introduction

Sentiment transfer [6] aims at altering the sentiment polarity of a text while preserving its vanilla content meanings, which has been widely employed in a myriad of applications such as news sentiment transformation [3], passage editing [15] and data augmentation [25]. For example, when inputting a text sequence “stale food and poor service”, the expected sentiment-transferred output would be “fresh food and good service”, which modifies the sentiment polarity from negative to positive while maintaining the content information.

Existing sentiment transfer models could be roughly categorized into two categories. The first type of work implicitly disentangles content and sentiment [8,17,3,2,13,9,1] by learning the latent representations of content and sentiment respectively, and then combining the content representation and target sentiment signal to generate transferred sentences. Such implicit approaches generally employ GAN (Generative Adversarial Network) to remove sentiment attributes from content representation and generate text indistinguishable from real data. However, previous works [7] noted that implicit methods generally suffer from inferior performance on content preservation and low interpretability. Another way of sentiment transfer is decoupling content and sentiment explicitly

[12,20,16,18,24,14,5,11], which first identifies sentiment-associated words and replaces them with words related to the target sentiment while keeping other words unchanged. Explicit methods benefit from the simplicity and explainability since the sentiment-tokens are explicitly uncovered and replaced.

Despite the promising performance of existing sentiment transfer models, they are still facing two crucial challenges. First, it is intractable to separate the sentiment style from the semantic content precisely. Existing implicit methods generally split the hidden representation of the input sequence into two vectors, which contain style and content information, respectively. However, such implicit methods generally perform unsatisfactorily on content preservation sub-task [7], probably brought by the loss of content information in the process of disentangling. Explicit models are capable of identifying emotion-associated tokens explicitly, while they generally locate these tokens simply based on prior linguistic knowledge and rules. Such heuristic methods are incapable of ensuring the precise correlations between tokens in masking positions and sentiment signals, leading to an obscure disentanglement. Second, it is nontrivial to combine content and target sentiment to generate a target sentence effectively. Existing works usually assume that the overall emotional tendency of the generated sentence should be close to the target sentiment label. [20] leverages Attribute Conditional Masked Language Model(AC-MLM) to fill in masked positions. [12] retrieves new phrases associated with target attributes from the corpus and combines them with content information. However, these methods only focus on sentence-level sentiment labels and pay no attention to the word-level polarities, thus cannot capture such fine-grained semantic information and the connections between sentence-level sentiment and word-level polarity.

In this paper, we propose a novel **Sentiment Transfer** model AM-ST to handle the mentioned challenges based on **Adaptive Masking**. First, we identify emotion-associated tokens and mask them using a trainable mask module, which is capable of adaptively learning the optimal mask positions. Then, a sentiment-aware masked language model is leveraged to fill in blanks in these masked positions, incorporating both context and sentiment polarity to improve transfer accuracy. Specifically, following the assumption of [20] that transferred sentence can be generated by simply replacing several emotional-related words, AM-ST adopts a mask classifier to identify the appropriate mask positions and then mask them with special tokens. In order to certify that sentiment information only appears in the masked tokens and not in the rest tokens, we design two types of losses: classification losses and adversarial losses, ensuring the clear separation of sentiment and content. After that, in the filling blanks phase, we adopt a sentiment-aware masked language model based on a reconstruction loss to predict both sentence- and token-level polarities in the masked positions. Experimental results on two popular datasets demonstrate the superiority of our proposal. Our major contributions are summarized as follows:

- To alleviate the challenge of unclear disentanglement brought by low identifying quality, we propose an adaptive masking module, setting mask position as a trainable parameter to certify accuracy in identifying sentiment words.

- We further propose a sentiment-aware masked language model in the infilling blanks stage, which captures both semantic context and emotional signals.
- Experimental results on popular datasets indicate that the proposed AM-ST model consistently outperforms SOTA models.

2 Related Work

Sentiment transfer is closely related to text style transfer, where the key is to modify the style of input text and retain the content. To address the problem of separating style and content, there are two major ways: 1) **Implicitly Separating:** These models divide the hidden representation of input text into two representations, style embedding and content embedding. Fu et al. [3] use adversarial networks to learn separated content representations and style representations. [16] generates pseudo-parallel corpora for text sentiment transfer. Yang et al. [23] replace the style discriminator in the adversarial learning with a language model. However, [10] experimentally proves that simply adopting GAN is inadequate to separate style from content, and it is easy to recover style information using content representation. 2) **Explicitly Separating:** This type of model explicitly identifies and replaces style-associated tokens. [12] removes original attribute phrases and retrieves new terms related to the target attribute. [20] converts the emotion transfer problem into a text fill-in-the-blank task through a pre-trained masked language model. [2] combines a pre-trained language model with attribute classifiers to guide text generation. Explicit methods benefit from their simplicity and explainability, because they clearly show which tokens are related to style. However, most current works exploit only prior language knowledge and rules to identify sentiment tokens without examining the association between masked tokens and sentiment information. Different from existing works, we adopt a trainable masked classification module to detect and mask the sentiment-related tokens.

3 Problem Definition

Given an input text sequence $X = \{x_1, x_2, \dots, x_N\}$ (x_i indicates a single word), its source sentiment label l , and the target sentiment label \tilde{l} , sentiment transfer aims to generate the target sentence \tilde{X} which maintains the content information of X while having the target sentiment label \tilde{l} . Following the previous work [8], we select the binary sentiment label set $\{positive, negative\}$.

4 Methodology

4.1 Framework

Figure 1 demonstrates the framework of the proposed AM-ST model. Given the input sentence *stale food and poor service*, AM-ST first utilizes the adaptive masking module to separate the sentiment words from other content words, resulting the masked sentiment token set S and the content token set C . After masking the sentiment tokens with the special tokens, the sentence is converted

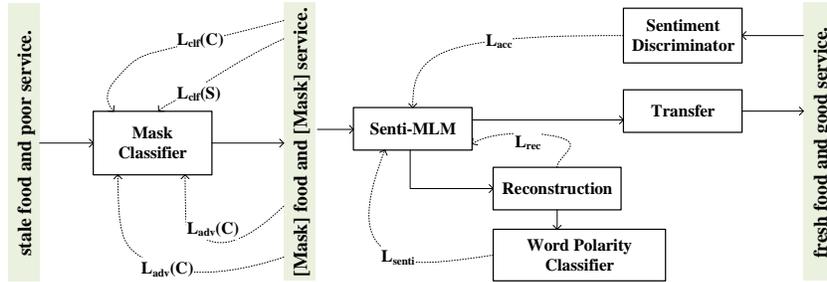


Fig. 1. Framework of the proposed AM-ST model.

to *[mask] food and [mask] service*. Then, based on the learned content text C and the target sentiment label \tilde{l} , we further propose a generative module to predict tokens to infill words conforming target sentiment. Both the word-level and sentence-level sentiment polarities are properly incorporated to generate desirable sentiment-transferred sentences. The input text is transferred into the positive sentiment “*fresh food and good service*”.

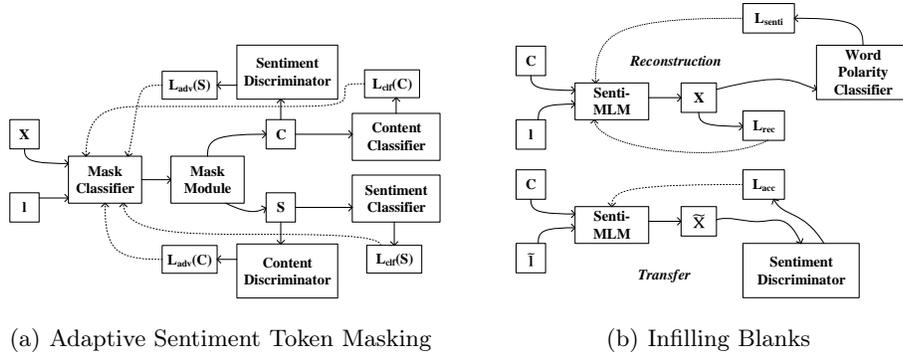
4.2 Adaptive Sentiment Token Masking

In this stage, we aim to identify and mask sentiment-associated tokens to achieve pure content text. Previous works generally rely on prior linguistic knowledge (e.g., sentiment dictionary, co-occurrence frequency of words with a particular sentiment) to identify the sentiment tokens. However, such rule-based heuristic methods might not be optimal. Downstream datasets might contain unique sentiment tokens, which might be inconsistent with the general knowledge and ignored by previous works, leading to low generality and inferior performance. Instead of directly leveraging prior linguistic knowledge to discover and mask sentiment tokens, we view the mask positions as the learnable variables. Our model is expected to automatically locate the positions of the sentiment words to better align with the domain-specific knowledge.

Given the input text sequence X and the vanilla sentiment label l , we first need to decide whether token x_i should be masked. A mask classifier as shown in Figure 2(a) is integrated to learn the masking probability of each token. The mask classifier is implemented as an attention-based classifier.

$$y = \text{softmax}(\alpha \cdot H) \quad (1)$$

where H denotes the hidden vectors generated by a Bi-LSTM model [4], and α denotes the attention weight vector generated based on the attention mechanism [19]. The output $y \in \mathbb{R}^{N \times 1}$ indicates the probability of x_i containing sentiment information and should be masked. If the predicted probability exceeds a threshold, we will mask the token in this position. We define S as the set of the tokens being masked, and C as the set of the rest of the tokens in X . The mask classifier will update its parameters adaptively during model training. In addition,



(a) Adaptive Sentiment Token Masking

(b) Infilling Blanks

Fig. 2. Two phases of the proposed AM-ST model.

we introduce extra constraints to facilitate disentangling between S and C . Intuitively, all tokens in S should contain sentiment information but no content information, while tokens in C should only include content information without any sentiment information. Such constraints are comprehensively satisfied by the following two types of losses.

Classification Loss. Classification losses ensure that the extracted token set C and S should capture the content and sentiment information, respectively. Thus, we design the following two classification losses for these two types of tokens.

Sentiment Classification Loss. The extracted set S should be associated with sentiment information, and we employ a pre-trained sentiment classifier $clf(S)$ implemented as a fully-connected layer with activation function Softmax to examine this. The input of the $clf(S)$ is S , and the loss function of $clf(S)$ is:

$$L_{clf}(S) = - \sum_{l \in \text{labels}} t_s(l) \log y_{s(S)}(l) \quad (2)$$

where $t_s \in \mathbb{R}^{N \times 2}$ denotes the output of $clf(S)$. $y_{s(S)}(l)$ denotes the probability of emotional polarity of S being l . $L_{clf}(S)$ is a cross-entropy loss and approximates the distance between the predicted distribution $y_{s(S)}$ and ground-truth distribution t_s . It is worth noting that classifier $clf(S)$ is pre-trained and its parameters will not be updated during the training of our model, so the reduction of $L_{clf}(S)$ is not caused by the classifier’s better predictive capacity, but brought by richer sentiment information encoded in the input text.

Content Classification Loss. To ensure the content information is captured by C , we design a content classification loss to measure the closeness between the vanilla input X and the content set C in terms of textual content. Following previous work [8], we use bag-of-words (BoW) to measure content completeness. To testify whether C preserves content well, we employ a Softmax content classifier $clf(C)$ to predict BoW distribution with the input of C . The loss function of $clf(C)$ is defined as the cross-entropy between ground-truth distribution t_c and predicted distribution $y_{c(C)}$:

$$L_{clf}(C) = - \sum_{w \in \text{vocabulary}} t_c(w) \log y_{c(C)}(w) \quad (3)$$

where $t_c(w) = \frac{\text{count}(w, X)}{N}$ is the ground-truth BoW distribution, $\text{count}(w, X)$ is the frequency of word w in the vocabulary appearing in X . $y_{c(C)}(w)$ denotes the predicted probability of word w 's appearance.

Adversarial Loss. To separate content and sentiment information, it is indispensable to examine whether C and S contain overlapping information. We further introduce two adversarial losses to accomplish the objective.

Sentiment Adversarial Loss. To ensure that content set C contains as little sentiment information as possible, we adopt a two-step training paradigm. In the first step, we introduce a Softmax sentiment discriminator $dis(S)$ to predict the sentiment label of C . To improve the performance of $dis(S)$, we introduce $L_{dis}(S)$ to measure the distance between predicted distribution $y_{s(C)}$ and ground-truth distribution t_s :

$$L_{dis}(S) = - \sum_{l \in \text{labels}} t_s(l) \log y_{s(C)}(l) \quad (4)$$

where $y_{s(C)}(l)$ represents the probability of the predicted label of C to be l . As $dis(S)$ has been trained to predict the sentiment label using C , in the second step, we introduce $L_{adv(S)}$ to punish the classification ability of $dis(S)$, which is implemented as the entropy of $y_{s(C)}$ and measures the accuracy of prediction. $L_{adv(S)}$ is maximized when $y_{s(C)}$ is evenly distributed, which means that sentiment labels are completely unpredictable for C :

$$L_{adv}(S) = - \sum_{l \in \text{labels}} y_{s(C)}(l) \log y_{s(C)}(l) \quad (5)$$

Content Adversarial Loss. In order to remove content information in S , we introduce $L_{adv}(C)$, which is calculated in the similar manner of $L_{adv}(S)$. First, we adopt a content Softmax discriminator $dis(C)$ and optimize its ability to predict content BoW contribution $y_{c(S)}$ using S by minimizing $L_{dis}(C)$.

$$L_{dis}(C) = - \sum_{w \in \text{vocabulary}} t_c(w) \log y_{c(S)}(w) \quad (6)$$

where t_c is the ground-truth distribution and $y_{c(S)}$ is the predicted distribution of $dis(C)$. Then, we punish the discernment of content discriminator. $L_{adv}(C)$ is the entropy of $y_{c(S)}$ and achieves its maximum when it is impossible to discern content distribution using S :

$$L_{adv}(C) = - \sum_{w \in \text{vocabulary}} y_{c(S)}(w) \log y_{c(S)}(w) \quad (7)$$

Overall objective function Based on the previous losses, the final objective function is formally designed as:

$$L_{total} = \lambda_1 L_{clf}(S) - \lambda_2 L_{adv}(S) + \lambda_3 L_{clf}(C) - \lambda_4 L_{adv}(C) \quad (8)$$

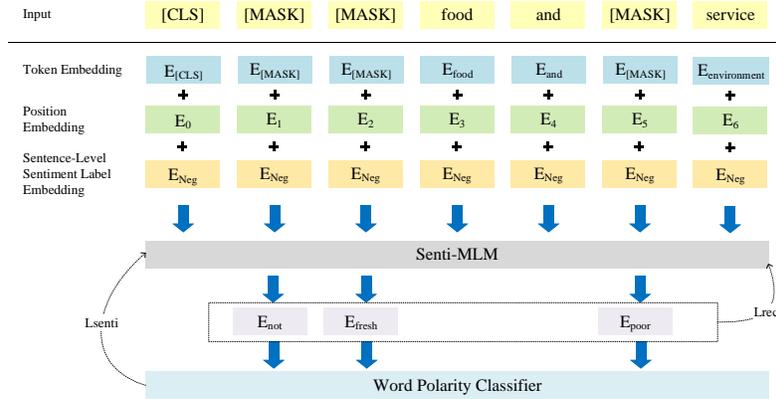


Fig. 3. Sentiment aware masked language model (Senti-MLM). Token embeddings, position embeddings and sentence-level sentiment embeddings are the inputs, and Senti-MLM is trained to predict tokens and word-level polarities in the masked positions.

Where $\lambda_1, \lambda_2, \lambda_3$ and λ_4 are the weights of the corresponding losses.

It is worth noting that although both the attention-based method in [20] and our model utilizes classifiers to find mask position, the former classifier is not updated synchronously during model training, but trained in advance. Therefore, the performance of [20] cannot be continuously improved as it ignores the feedback on whether masked tokens are related to sentiment. Instead, our mask position classifier adaptively updates its parameters using the above four losses and continues improving identifying accuracy.

4.3 Infilling Blanks

In this stage, our model will infill tokens in masked positions using a sentiment-aware masked language model (Senti-MLM) as shown in Figure 2(b). Although MLM performs well in clozing tasks, it only considers contextual information and overlooks sentiment information. However, antonyms with opposite sentiment polarities tend to have different contexts. [21] proved that MLM performs well in learning features of domains and semantics, but tends to neglect opinion words and sentiments. Therefore, simply utilizing MLM to infill words in masked positions may not achieve desirable performance since tokens in these positions still contain rich sentiment information. To solve this problem, we adopt Senti-MLM, which is able to predict the tokens in the masked positions considering not only context but also sentiment information to facilitate transfer accuracy.

On the basis of the MLM of the pre-trained BERT model, we design a new paradigm to incorporate sentiment information. As shown in Figure 3, the traditional segmentation embeddings are replaced by sentence-level sentiment label embeddings. The training objective is expected to predict tokens in masked positions and their word-level polarities. In the adaptive masking stage, we have proved that all masked tokens are associated with sentiment information. Therefore, the proposed Senti-MLM overcomes MLM’s shortcoming of overlooking

sentiment information and performs well in using words with proper sentiment polarity to fill in these blanks.

We employ a reconstruction task to train Senti-MLM. Given the set of content tokens C and the original label l , the objective aims to reconstruct the input sentence X using Senti-MLM. We measure the performance of the reconstruction task by L_{rec} , which reflects the content preservation ability of Senti-MLM:

$$L_{rec} = \sum_{t_i \in S} p(t_i | l; C) \quad (9)$$

To take word-level polarities into consideration, we train a word polarity classifier to predict word-level polarities using the hidden-state h generated by Senti-MLM, and its output is y_h . Then we introduce L_{senti} to measure Senti-MLM’s ability to recover word-level polarities by comparing y_h and ground-truth word-level polarity t_{polar} .

$$L_{senti} = - \sum_{l \in \text{labels}} t_{polar}(l) \log y_h(l) \quad (10)$$

Then, loss L_{rec} and L_{senti} are weighted combined to finetune Senti-MLM, and ϑ_1 and ϑ_2 denote the corresponding weights:

$$L_1 = \vartheta_1 L_{rec} + \vartheta_2 L_{senti} \quad (11)$$

After that, we leverage Senti-MLM to perform the transfer task. Content tokens C and target sentence label \tilde{l} are fed into Senti-MLM, and the transferred sentence \tilde{X} can be generated. To evaluate the transfer accuracy, we introduce L_{acc} to measure whether \tilde{X} conforms with target label \tilde{l} :

$$\tilde{X} = \text{SentiMLM}(C, \tilde{l}) \quad (12)$$

$$L_{acc} = - \log p(l | \tilde{X}) \quad (13)$$

We continue finetuning Senti-MLM using L_{rec} and L_{acc} . The overall loss in this stage is $L_2 = \vartheta_3 L_{rec} + \vartheta_4 L_{acc}$, where ϑ_3 and ϑ_4 trade off between L_{rec} and L_{acc} .

5 Experiment

5.1 Experimental Settings

Dataset Following previous works [12], we adopt two popular datasets, Yelp⁵ and Amazon⁶, to evaluate the performance of our proposal. Yelp contains business reviews in which each review is labeled with negative or positive sentiment. Similarly, Amazon dataset contains product reviews from Amazon, each of which is manually labeled as either negative or positive.

⁵ <https://www.yelp.com/dataset>

⁶ <https://www.kaggle.com/datasets/bittlingmayer/amazonreviews>

Table 1. Dataset statistics.

Dataset	Labels	Train	Valid	Test
Yelp	Positive	270K	2000	500
	Negative	180K	2000	500
Amazon	Positive	277K	985	500
	Negative	278K	1015	500

Baselines We compare the proposed AM-ST model with the following popular baselines for verifying the performance.

- CrossAligned [17]: CrossAligned generates the original sentence back using a generator that combines content representation with the original label.
- StyleEmbedding [3]: Style embedding is fed into a decoder to generate text given different target sentiment signals.
- MultiDecoder [3]: MultiDecoder is a seq2seq model using multiple decoders. Each decoder independently generates a corresponding text style.
- CycledReinforce [22]: This model consists of a deemotionalizing module and an emotionalizing module, which extracts non-emotive semantic information and then emotionalizes neutral sentences.
- DeleteAndRetrieval [12]: DeleteAndRetrieval removes original attribute phrases and retrieves new terms related to the target attribute in the corpus.
- DisentangledRepresentation [8]: VAE with auxiliary multitask and adversarial objectives are used to learn content embeddings and style embeddings.
- AC-MLM-Frequency [20]: AC-MLM-Frequency converts the emotion transfer problem into a clozing task through a masked language model.
- AC-MLM-Fusion [20]: An extension of AC-MLM-Frequency which employs an attention mechanism to further filter retrieved sentiment words.

Evaluation Metrics Following the previous work [12], we select Accuracy and BLEU as the evaluation metrics. Accuracy is calculated by how likely a transferred sentence conforms with the target sentiment, which is an indicator of transfer accuracy. BLEU is computed by the similarity between human reference by [12] and the generated transferred sentence. A high BLEU score indicates that the model performs well in content preservation.

5.2 Implementation Details

In the adaptive sentiment token masking phase, we use an attention classifier pre-trained on Yelp and Amazon datasets as the mask classifier. $clf(S)$, $clf(C)$, $dis(S)$ and $dis(C)$ are implemented using four Softmax classifiers. In the phase of infilling blanks, we use the pre-trained Bert_{base} as the checkpoint. The word polarity classifier is a fully-connected layer with a 768×3 tensor as input, and the sentiment discriminator is implemented as a convolutional neural network (CNN). We first finetune the checkpoint to fit Senti-MLM as well as train the word polarity classifier for ten epochs, only L_{rec} is computed in this stage. Then we further train six epochs to finetune Senti-MLM, where both L_{rec} and L_{senti}

Table 2. Experimental results on Amazon and Yelp.

	Yelp		Amazon	
	ACC(%)	BLEU	ACC(%)	BLEU
CrossAligned [17]	73.1	3.1	74.1	0.4
StyleEmbedding [3]	8.7	11.8	43.3	10.0
MultiDecoder [3]	47.6	7.1	68.3	5.0
CycledReinforce [22]	85.2	9.9	77.3	0.1
DeleteAndRetrieval [12]	88.7	8.4	48.0	22.8
DisentangledRepresentation [8]	91.5	12.2	82.4	25.2
AC-MLM-Frequency [20]	95.1	11.6	64.5	27.2
AC-MLM-Fusion [20]	95.3	12.3	85.2	28.3
AM-ST	97.1	12.9	86.4	29.7

constitute the total loss in this stage. Finally, we finetune Senti-MLM and train the sentiment discriminator for another ten epochs, where L_{rec} and L_{acc} are used to update the parameters of Senti-MLM and sentiment discriminator.

As shown in Equation (8), λ_1 , λ_2 , λ_3 and λ_4 control the trade-off among four losses in adaptive masking stage. Similarly, hyper-parameters ϑ_1 , ϑ_2 , ϑ_3 and ϑ_4 in Equation (11) also define the weights of different losses. We carefully tune these hyperparameters on the validation set, and report the testing results of the parameter setting with the best validation performance. The tuned value of λ_1 , λ_2 , λ_3 and λ_4 are 0.2, 0.1, 0.4 and 0.3, respectively. Best validation performance is achieved when ϑ_1 , ϑ_2 , ϑ_3 and ϑ_4 are 0.4, 0.2, 0.1 and 0.3.

5.3 Quantitative Analysis

All sentiment transfer models are evaluated five times, and the average performance is reported in Table 2. For the StyleEmbedding model, the parameters of the encoder and decoder are fixed when generating sentences of target sentiments, leading to inferior sentiment capability and thus achieving the worst performance. AC-MLMs are the strongest baselines because they take advantage of MLM’s strong ability to predict masked tokens according to the semantic context. AC-MLM-Fusion achieves better performance since it further overcomes the deficiency of AC-MLM-Frequency by introducing an attention-based classifier to filter pseudo-sentiment words. One can clearly see that our proposal consistently outperforms baseline models on both datasets, verifying the effectiveness of the proposed adaptive masking mechanism and the token-level sentiment polarity.

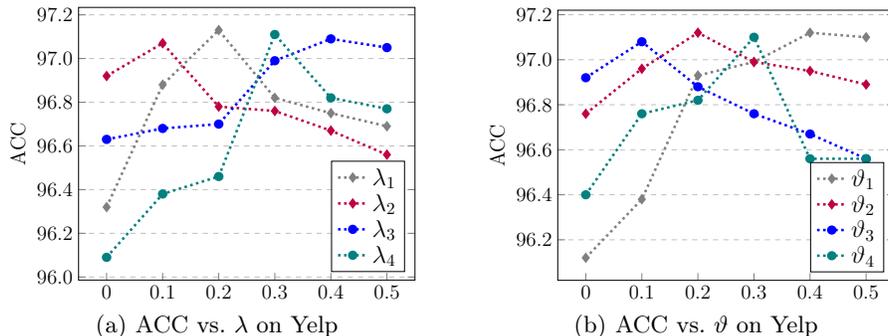
5.4 Ablation Study

We further conduct an ablation study to verify the effectiveness of different components. Five ablation models are designed by removing different objective functions, namely sentiment classification loss in Equation (2), content classification loss in Equation (3), sentiment adversarial loss in Equation (5) and content adversarial loss in Equation (7). Table 3 presents the experimental results of different ablation models. One can easily observe that model performance significantly decreases after removing any components, verifying these modules are

Table 3. Ablation Study.

	Yelp		Amazon	
	ACC(%)	BLEU	ACC(%)	BLEU
$-L_{clf}(S)$	96.1	12.7	84.7	29.0
$-L_{clf}(C)$	96.3	11.7	85.9	28.6
$-L_{adv}(S)$	96.4	12.4	85.8	29.4
$-L_{adv}(C)$	96.8	12.1	86.2	28.9
$-L_{senti}$	95.6	12.2	85.2	29.2
AM-ST	97.1	12.9	86.4	29.7

indispensable to a successful sentiment transfer. It is reasonable as the classification losses contribute to capturing the sentiment/content information, while the adversarial losses are capable of separating the sentiment and content information. The last module L_{senti} incorporates the token-level sentiment signals into the MLM process, facilitating the generation phase.

**Fig. 4.** Parameter sensitivity analysis on eight core hyper-parameters.

5.5 Parameter Sensitivity Analysis

Here we study the performance sensitivity of our proposal on eight core parameters: the weights λ_i in Equation (8) and the weights ϑ_i of Formula (11). As the performance trends on the two datasets are similar, here we only report the results on the Yelp dataset. We first fix other hyper-parameters and then report the results by tuning the target hyper-parameter in the range of $[0, 0.5]$. Figure 4(a) and 4(b) presents the experimental results. One can see that with the increase of different hyper-parameters, the performance over all datasets first increases and then decreases, leading to a similar tendency. Thus, these hyper-parameters should be carefully tuned to achieve desirable performance.

6 Conclusion

We present a novel model for sentiment transfer, which views the mask positions as trainable parameters to accurately identify and mask sentiment-related words. In addition, a sentiment-aware masked language model is adopted to infill blanks more efficiently by considering both context and word-level polarity. Experiments demonstrate that our model consistently outperforms SOTA models.

References

1. Chen, L., Dai, S., Tao, C., Zhang, H., Gan, Z., Shen, D., Zhang, Y., Wang, G., Zhang, R., Carin, L.: Adversarial text generation via feature-mover’s distance. *Advances in Neural Information Processing Systems* **31** (2018)
2. Dathathri, S., Madotto, A., Lan, J., Hung, J., Frank, E., Molino, P., Yosinski, J., Liu, R.: Plug and play language models: A simple approach to controlled text generation. In: *ICLR* (2020)
3. Fu, Z., Tan, X., Peng, N., Zhao, D., Yan, R.: Style transfer in text: Exploration and evaluation. In: *AAAI* (2018)
4. Graves, A., Schmidhuber, J.: Framewise phoneme classification with bidirectional lstm networks. In: *IJCNN*. pp. 2047–2052 (2005)
5. Guu, K., Hashimoto, T.B., Oren, Y., Liang, P.: Generating sentences by editing prototypes. *Transactions of the Association for Computational Linguistics* **6**, 437–450 (2018)
6. He, R., McAuley, J.: Ups and downs: Modeling the visual evolution of fashion trends with one-class collaborative filtering. In: *WWW*. pp. 507–517 (2016)
7. Hu, Z., Lee, R.K.W., Aggarwal, C.C., Zhang, A.: Text style transfer: A review and experimental evaluation. *ACM SIGKDD Explorations Newsletter* pp. 14–45 (2022)
8. John, V., Mou, L., Bahuleyan, H., Vechtomova, O.: Disentangled representation learning for non-parallel text style transfer. In: *ACL*. pp. 424–434 (2019)
9. Krishna, K., Nathani, D., Samanta, B., Talukdar, P.: Few-shot controllable style transfer for low-resource multilingual settings. In: *ACL*. pp. 7439–7468 (2022)
10. Lample, G., Subramanian, S., Smith, E., Denoyer, L., Ranzato, M., Boureau, Y.L.: Multiple-attribute text rewriting. In: *ICLR* (2018)
11. Lee, J.: Stable style transformer: Delete and generate approach with encoder-decoder for text style transfer. In: *Proceedings of the 13th International Conference on Natural Language Generation*. pp. 195–204 (2020)
12. Li, J., Jia, R., He, H., Liang, P.: Delete, retrieve, generate: a simple approach to sentiment and style transfer. In: *NAACL*. pp. 1865–1874 (2018)
13. Liu, A., Wang, A., Okazaki, N.: Semi-supervised formality style transfer with consistency training. In: *ACL*. pp. 4689–4701 (2022)
14. Madaan, A., Setlur, A., Parekh, T., Poczós, B., Neubig, G., Yang, Y., Salakhutdinov, R., Black, A.W., Prabhume, S.: Politeness transfer: A tag and generate approach. In: *ACL*. pp. 1869–1881 (2020)
15. Parra G, L., Calero S, X.: Automated writing evaluation tools in the improvement of the writing skill. *International Journal of Instruction* **12**(2), 209–226 (2019)
16. Shang, M., Li, P., Fu, Z., Bing, L., Zhao, D., Shi, S., Yan, R.: Semi-supervised text style transfer: Cross projection in latent space. In: *EMNLP*. pp. 4937–4946 (2019)
17. Shen, T., Lei, T., Barzilay, R., Jaakkola, T.: Style transfer from non-parallel text by cross-alignment. *NeurIPS* (2017)
18. Sudhakar, A., Upadhyay, B., Maheswaran, A.: “transforming” delete, retrieve, generate approach for controlled text style transfer. In: *EMNLP*. pp. 3269–3279 (2019)
19. Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A.N., Kaiser, Ł., Polosukhin, I.: Attention is all you need. *NeurIPS* (2017)
20. Wu, X., Zhang, T., Zang, L., Han, J., Hu, S.: Mask and infill: Applying masked language model for sentiment transfer. In: *IJCAI*. pp. 5271–5277 (2019)
21. Xu, H., Shu, L., Yu, P., Liu, B.: Understanding pre-trained BERT for aspect-based sentiment analysis. In: *COLING*. pp. 244–250 (2020)

22. Xu, J., Sun, X., Zeng, Q., Zhang, X., Ren, X., Wang, H., Li, W.: Unpaired sentiment-to-sentiment translation: A cycled reinforcement learning approach. In: ACL. pp. 979–988 (2018)
23. Yang, Z., Hu, Z., Dyer, C., Xing, E.P., Berg-Kirkpatrick, T.: Unsupervised text style transfer using language models as discriminators. NeurIPS (2018)
24. Zhang, Y., Xu, J., Yang, P., Sun, X.: Learning sentiment memories for sentiment modification without parallel data. In: EMNLP. pp. 1103–1108 (2018)
25. Zheng, X., Chalasani, T., Ghosal, K., Lutz, S., Smolic, A.: Stada: Style transfer as data augmentation. In: VISAPP (2019)