

1 **Dear Reviewers,**

2 Thanks a lot for your helpful and insightful comments. We read them carefully and tried our best to address the issues
3 raised in the comments. Below are our responses to the key issues in the reviews:

4 **To Reviewer #1:**

5 **1.** (For the problem that "How to ensure the classifier has the optimal set of parameters as z is changed") The classifier's
6 parameters are kept fixed as z is changed. We train the autoencoder and the classifier on the training set, which is
7 diverse and contains texts of varying degrees of attributes, reflected by the different confidence values given by the
8 classifier. Thus we believe when the text in the test set is encoded into the same latent space as the training set, the
9 classifier can continue to provide the optimal direction of modification, even when z changes.

10 **2.** (For the problem that "It's unclear what the notion 'degree of transfer' means") The 'degree of transfer' here means
11 the apparent extent of the desired attributes embodied in the transfer results. Taking sentiment transfer as example, it
12 is reflected in the sentimental intensity felt from the text (e.g., weak positive, moderate positive and strong positive).
13 Different from most previous work that only provides binary control over attributes, one advantage of our model is
14 the ability to give control over the degree of attribute transfer desired. Following the evaluation settings of previous
15 work, we use various automatic evaluation indicators (ie, Acc, BLEU, and PPL in Section 4.1) to automatically evaluate
16 transfer results from different aspects. Particularly, 'Acc' is used to evaluate the attribute's accuracy. In Figure 2,
17 by showing the values of three automatic indicators under different fixed modification weights, we demonstrate the
18 influence of the weight on controlling the transfer results. Sorry for the misunderstanding, we will add more words to
19 explain this.

20 **3.** Thanks for your suggestion, we will add details about the interaction between the autoencoder and the classifier into
21 Algorithm 1.

22 **4.** (For the problem that "What does 'this may undermine the integrity' mean ?") The word "integrity" here means
23 naturality. For example, some phrase-based methods directly delete/replace/insert some sentimental words/phrases,
24 which may result in the generated sentences that are unnatural and not as likely to be written by the human (can be
25 partly reflected by the automatic indicator 'PPL'). Sorry for the misunderstanding, we will add more explanations.

26 **5.** Thank you for your careful review, we will further verify our model on more involved Multi-Aspect Sentiment
27 Transfer dataset and multi-field conditioned generation dataset by Ficler and Goldberg 2017. We will improve our paper
28 by adding more explanations about "degree of transfer" and "Automatic Evaluation" results, adding more comprehensive
29 references, and carefully proofreading the paper.

30 **To Reviewer #2:**

31 **1.** (For concerns #1 and #2) We will compare our model with some simpler autoencoder architectures (e.g., LSTM,
32 GRU), and show the BLEU scores between input and reconstructions. We will add more ablation&comparison
33 experiments, to show the gains of different parts of our model.

34 **2.** (For concerns #3) Thanks for your insightful suggestions, we will add references about activation maximization
35 including PPGN, DeVise, etc., and provide a more comprehensive related work section.

36 **To Reviewer #3:**

37 **1.** (For the problem that "how human evaluation was performed") For each test sample, we hired 3 workers to annotate
38 and average the scores. Moreover, we divided the test samples of each model on each task into the same number of
39 small sets, and the same person annotated the same task for all the models. Due to a large number of samples to be
40 labeled, we did not perform an A/B test, and we will add it if necessary. There was some test data collected for acquire
41 human references for attribute transfer, and the references written by human for BLEU calculation are provided by
42 previous work (Juncen Li, Robin Jia, He He, and Percy Liang. Delete, retrieve, generate: a simple approach to sentiment
43 and style transfer.). Thanks for your suggestions, we will provide a clearer explanation to explain this.

44 **2.** (For the problem that "The argument about the correlation between human evaluation and automatic evaluation") The
45 results of human evaluation and automatic evaluation are not always consistent with each other. In this study, we follow
46 the evaluation settings of previous works to provide comprehensive results with various indicators. Thank you for your
47 careful review, we will remove this argument about the correlation but provide more detailed evidence and analysis.