

## A Appendix

### A.1 Implementation Details

We discuss implementation details including model parameters. In the graph encoder, the maximum number of nodes is set as 64 for all input graphs. The dimension is set as 100 in forward hidden state vectors  $\mathbf{h}_{k,i}^{\text{fwd}}$  and backward hidden state  $\mathbf{h}_{k,i}^{\text{bwd}}$  for any node  $i$  in any iteration  $k$ . The total number of iterations  $K$  is set as 6 in our experiments. For the sequence decoder, the dimension of hidden state vectors is 400 in its recurrent unit.

As for the discriminators, we use the default parameters set in the baseline. The style discriminator has 50-dimension hidden state vectors and 20-dimension attention vectors. The language model serving as the fluency discriminator also has 50-dimension hidden state vectors.

We pre-train the model for 4 epochs with a batch size of 16 and a learning rate of  $1e - 4$ , and train it for 2 epochs with learning rate of  $1e - 5$  in reinforcement learning. During reinforcement learning, the sample size of target sentences is 4 when the reward  $Q$  is estimated from the evaluation scores  $r$  of complete target sentences. One epoch takes about 2 hours in pre-training, and takes about 5 hours in reinforcement learning on one GPU.

### A.2 More Examples of Transferred Sentences

Table 6 shows some transferred sentences of four models—CA, MD, RL and GT—for both positive-to-negative and negative-to-positive transfer.

### A.3 Human annotations

Human annotations are complementary evaluation to automatic metrics. In terms of the semantic preservation, the automatic evaluation tends to compare the lexical overlap between input sentences and transferred sentences. Therefore, it favors the transferred outputs which contains the same words as the inputs. Table 6 shows an example with the original sentence “keep up the great work!”. CA and RL receive higher score than GT in semantic preservation given that they retain many words from the input, while GT gives a better transferred sentence. Another example in Table 6 is the one with original sentence “i love everything about this place”. Both RL and GT expresses negative feelings about a place, but the automatic metric assigns a much higher score to RL.

As for the automatic evaluation of transfer strength, the metric tends to be affected by senti-

	Semantic	Style	Fluency
CA wins	0.19	0.16	0.20
MD wins	0.09	0.21	0.15
RL wins	0.15	0.17	0.19
GT wins	<b>0.38</b>	<b>0.26</b>	<b>0.23</b>
Tie	0.19	0.20	0.23

Table 5: Percentage of model wins and ties.

ment words in transferred sentences. A negative-to-positive example in Table 6 has an input sentence “the food tasted sub-par at best”. The outputs of CA and RL contain the phrase “at best” which does not express the positive sentiment. However, their outputs are assigned with good style scores since the pre-trained style classifier considers “best” as a positive signal. While the automatic metrics are scalable to large datasets, they have limitations in the accuracy of evaluation. Hence human evaluation is adopted in complement to automatic metrics.

Based on the overall human evaluation of both negative-to-positive and positive-to-negative transfer, we report the percentage of wins for each model and the ties between multiple models in Table 5. As can be seen, GT outperforms other baselines in semantic preservation, style strength as well as fluency.

Transfer	Sentences
Negative-to-Positive	<p>Orig: overall , this place is really lax , horribly managed and not very clean .  CA: but , this place is very clean , and always clean .  MD: great bar and i recommend best and will be back of best .  RL: overall , this place is really clean , and really nice  GT: overall ... this place is really packed , its very clean .</p>
	<p>Orig: the food tasted sub-par at best .  CA: the food at best bbq .  MD: this is the best and will be back .  RL: the food tasted delicious at best .  GT: the food tasted fresh !</p>
	<p>Source: no attention to customers .  CA: great to work .  MD: overall , i will be back !  RL: great job to help .  GT: absolutely adore their customers .</p>
	<p>Source: associate was an ass with an attitude .  CA: staff was a amazing attitude .  MD: but it was so my favorite place to the best of i had .  RL: this was a fantastic job  GT: customer service is always on top of all levels .</p>
Positive-to-Negative	<p>Orig: yummy pizza and really good service .  CA: gross pizza but really bad and terrible .  MD: they are , it is , the food .  RL: gross pizza and really really bad service  GT: horrible service , terrible food .</p>
	<p>Orig: will definitely come back when i am in town .  CA: will not come back if i am in town in town .  MD: overall , do n't be worth it like _num_ .  RL: will not go back , i am in town in town  GT: do n't waste your time or money here .</p>
	<p>Orig: i love everything about this place .  CA: i just nothing about that about this place .  MD: she just the worst customer service i wo n't be .  RL: i hate everything about about this place  GT: i hate this place .</p>
	<p>Orig: keep up the great work !  CA: keep up the whole work work .  MD: is the staff is the food , i did n't even to me .  RL: keep up the work work work !  GT: the service sucked .</p>

Table 6: Examples of transferred sentences.