

# TARA: Token-level Attribute Relation Adaptation for Multi-Attribute Controllable Text Generation

Anonymous ACL submission

## Abstract

Multi-attribute controllable text generation (CTG) aims to generate fluent text satisfying multiple attributes, which is an important and challenging task. The majority of previous research on multi-attribute ignored the attribute relations. Recently, several work considers the attribute relations by explicitly defining them as "prohibitory". We argue that the attribute relations are not fixed and manifested as both "prohibitory" and promotive. To enhance multi-attribute CTG, in this paper, we propose TARA, which tackles multi-attribute controllable text generation with token-level attribute relation adaptation and representation, and uses a dynamic text generation strategy to exploit multi-attribute relations with balanced attribute control. We also define token-level attribute representation for multi-attribute CTG. Experimental results show that TARA achieves competitive control ability and comparable text quality and diversity over baseline methods.

## 1 Introduction

Multi-attribute controllable text generation (CTG) aims to generate fluent text satisfying multiple attributes. Previously, the majority of research on multi-attribute CTG didn't explicitly consider the relation between the attributes, which is a fundamental issue in multi-attribute CTG.

Recent work (Qian et al., 2022; Ding et al., 2023) take the attribute relations into consider and utilizes prefix tuning or VAE to train a multi-attribute model. Several work (Gu et al., 2022; Huang et al., 2023) further defines multi-attribute relation as inhibitory. For instance, Dist. Lens (Gu et al., 2022) identify that mutual interference of controllers causes attribute degeneration and searches for intersections in the attribute space. Prompt-Gating (Huang et al., 2023) use trainable gates to normalize the interference among attributes.

In fact, the attribute relations are not fixed, nor are they only manifested as "prohibitory". Take the

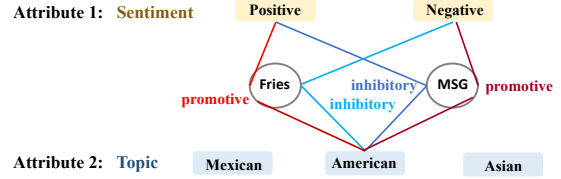


Figure 1: Token-level multi-attribute promotive and inhibitory relations. <sup>1</sup>

examples in Figure 1, MSG (Monosodium Glutamate) demonstrates the promotive relation between *negative* and *American*, and the inhibitory relation between *positive* and *American*. Therefore, multi-attribute CTG needs token-level attribute relation adaptation and exploitation of multi-attribute relation.

In this paper, we tackle the multi-attribute CTG with Token-level Attribute Relation Adaptation and representation, and propose TARA, which uses a dynamic text generation strategy. In summary, our contributions are as follows:

- We firstly identify both promotive and inhibitory attribute relations, and develop a token-level attribute relation adaptation method for multi-attribute CTG.
- The proposed attribute-adaptive prefix tuning adjusts attribute's expression with token-level attribute representation, and the dynamic text generation strategy we design balances multi-attribute control with promotive and inhibitory attribute relations.
- Experimental results verify that TARA performs better than existing methods on control ability and achieves text quality and diversity performances comparable with existing methods.

<sup>1</sup>MSG is the abbreviation of Monosodium Glutamate.

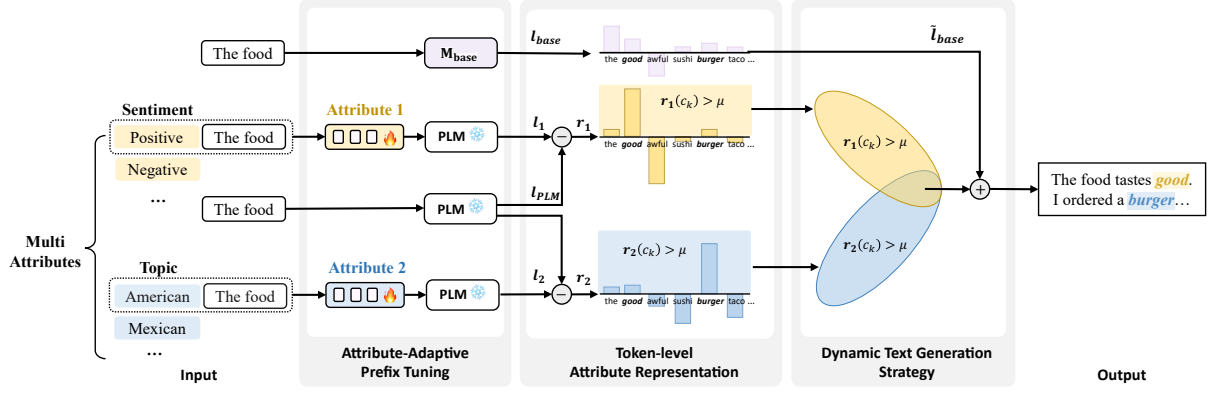


Figure 2: Overview of Multi-Attribute Relation Adaptation Method for CTG.

## 2 Proposed Method

We propose a novel multi-attribute relation adaptation method TARA for controllable text generation. Figure 2 illustrated the overall structure of TARA, which consists of attribute-adaptive prefix tuning, token-level attribute representation and dynamic text generation strategy.

### 2.1 Attribute-Adaptive Prefix Tuning

In multi-attribute CTG, the presence of tokens that exhibit inhibitory relations between attributes can lead to a degradation in the performance of multi-attribute control. Therefore, we employ attribute-adaptive prefix tuning for attribute models to weaken the inhibitory relations. This approach reduces the attribute expression of tokens with weak attribute expression, therefore enhancing the effectiveness of multi-attribute control.

Given the prompt or the current sequence as  $S_1^{t-1}$ , at the current time step  $t$ , we can obtain the attribute model’s logits  $l_t$  over the vocabulary  $\mathcal{V}$  through the language model. The language model will generate next token  $s_t$  by sampling  $s_t \sim P(s_t | S_1^{t-1})$  based on its logits  $l_t$ . We first convert the logits  $l_t$  to probabilities  $P(s_t | S_1^{t-1})$  with a temperature coefficient (see Equation 1) to make the attribute distinctions among tokens more apparent.

$$P(s_t | S_1^{t-1}) = \text{softmax}\left(\frac{l_t}{\tau}\right) \quad (1)$$

where  $\tau$  is the temperature coefficient.

Then, we use the L2 norm to constrain the probability distribution of the attribute model, avoiding extreme probability values and enhancing stability. Finally, we use the negative variance to enhance the attribute model’s ability to distinguish tokens with

varying degrees of attribute expression, mitigating tokens with weak single attribute expression. The adapt loss  $\mathcal{L}_{\text{adapt}}$  is defined as:

$$\mathcal{L}_{\text{adapt}} = \left( \|P(s_t | S_1^{t-1})\|_2 - \lambda_{\text{var}} \cdot \text{Var}(P(s_t | S_1^{t-1})) \right) \quad (2)$$

where  $\|P(s_t | S_1^{t-1})\|_2$  denotes the L2 norm of the probabilities, and  $\text{Var}(P(s_t | S_1^{t-1}))$  represents the variance of the probabilities. The term  $\lambda_{\text{var}}$  is the regularization weights for the L2 norm and the variance, respectively.

Finally, the total loss  $\mathcal{L}_{\text{total}}$  is defined as:

$$\mathcal{L}_{\text{total}} = -\log P(s_t | S_1^{t-1}) + \lambda_{\text{reg}} \mathcal{L}_{\text{adapt}} \quad (3)$$

where the terms  $\lambda_{\text{reg}}$  is the regularization weights for  $\mathcal{L}_{\text{adapt}}$ .

### 2.2 Token-level Attribute Representation

In TARA, we adapt fine-tune the corresponding attribute model for each attribute value using the same pre-trained language model. Given the current input  $S_1^{t-1}$ , we can obtain the logits distribution of the attribute model. Not only tokens that express the attribute characteristics will have high logits, but tokens that ensure text quality will also receive high logits, such as *with, a, the, is*. In TARA, we aim to utilize pure logits that only express the attribute characteristics for attribute control. Therefore, we define the attribute representation  $r_{\text{att}}$  as follows:

$$r_{\text{att}} = l_{\text{att}} - l_{\text{PLM}} \quad (4)$$

where att represents the attribute value,  $l_{\text{att}}$  represents the logits of attribute model,  $l_{\text{PLM}}$  represents the logits of base PLM. Attribute Representation

reflects the significance of the corresponding token in expressing the current attribute.

Suppose we have two attribute value  $i$  and  $j$  (TARA is capable of controlling more attributes). In the vocabulary  $\mathcal{V}$ , we define the tokens as  $\{c_1, \dots, c_n\}$ , where  $n$  is the total number of tokens in the vocabulary. We use  $c_k$  to represent a specific token in the vocabulary, where  $k \in \{1, \dots, n\}$ . The same token may demonstrate different attribute relations under different multi-attribute control. Therefore, we firstly set a threshold  $\mu$  to distinguish the different attribute expressions. Then we define two kinds of attribute value  $i$  and  $j$  relations at the token-level as follows:

**Promotive relation:**

$$\begin{cases} r_i(c_k) > \mu & \text{and} & r_j(c_k) > \mu \\ r_i(c_k) < \mu & \text{and} & r_j(c_k) < \mu \end{cases} \quad (5)$$

This indicates that the representations of  $r_i(c_k)$  and  $r_j(c_k)$  are consistent.

**Inhibitory relation:**

$$\begin{cases} r_i(c_k) > \mu & \text{and} & r_j(c_k) < \mu \\ r_i(c_k) < \mu & \text{and} & r_j(c_k) > \mu \end{cases} \quad (6)$$

This indicates that the representations of  $r_i(c_k)$  and  $r_j(c_k)$  are inconsistent.

In multi-attribute CTG, we establish dynamic vocabulary  $\mathcal{V}_i$  and  $\mathcal{V}_j$  for each attribute value. Then we can get:

$$\mathcal{V}_i = \{c_k \in \mathcal{V} \mid r_i(c_k) > \mu\} \quad (7)$$

$$\mathcal{V}_j = \{c_k \in \mathcal{V} \mid r_j(c_k) > \mu\} \quad (8)$$

## 2.3 Dynamic Text Generation Strategy

In MARA, we design a dynamic text generation strategy to exploit multi-attribute relation and balance multi-attribute control to steer the generation. We employ multi-attribute representations  $r_{\text{att}}$  in conjunction with the base logits  $l_{\text{base}}$  from a quality control model, which can either be an LLM or a small LM, sharing the same vocabulary as the attribute pre-trained model. In this setup, the attribute representations  $r_{\text{att}}$  manage multi-attribute control, while  $l_{\text{base}}$  ensures the quality of the text text. Besides, we design a dynamic weights to effectively balance the relations between two attributes, while ensure the text quality. Following Liu et al. (2021), we applied nucleus sampling

(Holtzman et al., 2020) to base model to obtain a fluent output sequence. At time step  $t$ , let  $\mathcal{V}' \subseteq \mathcal{V}$  represent the tokens included in the top- $p$  vocabulary of the base quality control model. The truncated logits  $\tilde{l}_{\text{base}}$  are

$$\tilde{l}_{\text{base}}[v] = \begin{cases} l_{\text{base}}[v] & \text{if } v \in \mathcal{V}' \\ -\infty & \text{otherwise} \end{cases} \quad (9)$$

For the two attribute value  $i$  and  $j$ , we define  $W_i$  and  $W_j$  represent the conditional probabilities of common tokens under attribute values  $i$  and  $j$ :

$$W_i = \text{softmax} \left( \frac{|\mathcal{V}_i \cap \mathcal{V}_j|}{|\mathcal{V}_i|} \right) \quad (10)$$

$$W_j = \text{softmax} \left( \frac{|\mathcal{V}_i \cap \mathcal{V}_j|}{|\mathcal{V}_j|} \right) \quad (11)$$

To normalize to multi-attribute weight to  $[0, 1]$  at the token-level, we design dynamic weights  $\tilde{W}_i$  and  $\tilde{W}_j$  as follows:

$$\tilde{W}_i = \frac{W_i}{W_i + W_j} \quad (12)$$

$$\tilde{W}_j = \frac{W_j}{W_i + W_j} \quad (13)$$

For a more reasonable sampling process, as in (Fan et al., 2018), we applied top- $K$  processing to the ensemble logits  $\tilde{l}_t$  during sample process. Therefore, the next token  $s_t$  can be obtained through the following dynamic text generation strategy:

$$\tilde{l}_t = \tilde{l}_{\text{base}} + (1 + \tilde{W}_i) \cdot r_i + (1 + \tilde{W}_j) \cdot r_j \quad (14)$$

$$\tilde{P}(s_t | S_1^{t-1}) = \text{softmax}(\tilde{l}_t), \quad (15)$$

$$s_t \sim \tilde{P}(s_t | S_1^{t-1}). \quad (16)$$

## 3 Experiments and Results

### 3.1 Experimental Setup

**Dataset** We choose widely used benchmark dataset YELP (Lample et al., 2019) for our experiments. Following previous work, we use sentiment attribute (positive and negative) and topic attribute (Asian, American and Mexican) for multi-attribute controllable text generation. Due to the page limit, please refer to Appendix B for more details about the experiment setup.

Method	Correctness (%)			Text Quality	Diversity
	Sentiment $\uparrow$	Topic $\uparrow$	Avg $\uparrow$	PPL $\downarrow$	mean-Dist $\uparrow$
PromptTuning* (Lester et al., 2021)	48.29	48.11	48.20	40.89	0.42
PrefixTuning* (Li and Liang, 2021)	47.53	69.11	58.32	147.47	0.31
ControlPrefixTuning (Clive et al., 2022)	58.98	45.36	52.17	89.80	0.48
GeDi* (Krause et al., 2021)	<b>99.47</b>	51.36	75.41	616.92	<b>0.75</b>
Tailor* (Yang et al., 2023)	80.68	68.72	74.70	40.29	0.39
Dist. Lens* (Gu et al., 2022)	77.47	66.98	72.22	52.59	0.26
PromptGating* (Huang et al., 2023)	84.80	<u>75.02</u>	<u>79.91</u>	<b>21.77</b>	0.42
<b>TARA-Multi (Ours)</b>	<u>90.17</u>	<b>80.32</b>	<b>85.25</b>	<u>40.16</u>	0.45

Table 1: The main results of multi-attribute CTG. For each method, we select 6 combinations (two sentiment attributes  $\times$  three topic attributes) as the final results. <sup>2</sup>

Variant	Correctness (%)			Text Quality	Diversity
	Sentiment $\uparrow$	Topic $\uparrow$	Avg $\uparrow$	PPL $\downarrow$	mean-Dist $\uparrow$
<b>TARA-Multi</b>	90.17	80.32	85.25	40.16	0.45
– Attribute-Adaptive Prefix Tuning	89.02	77.47	83.25	44.17	0.45
– Dynamic Text Generation Strategy	87.63	75.11	81.37	31.80	0.43

Table 2: Ablation Study of attribute-adaptive prefix tuning and dynamic text generation strategy of TARA.

**Evaluation Metrics** Following Yang et al. (2023); Huang et al. (2023), we conduct automatic and human evaluations for controllable accuracy and text quality. We conduct automatic evaluation from three aspects: (1) **Correctness** We finetune a sentiment classifier and topic classifier based on RoBERTa (Liu et al., 2019) for the evaluation of sentiment and topic accuracy. (2) **Text Quality** We calculate the perplexity (PPL) using GPT-2<sub>medium</sub> (Radford et al., 2019) to evaluate the fluency. (3) **Text Diversity** We use averaged distinctness (Li et al., 2015) to evaluate the diversity. We conduct human evaluation for sentiment relevance, topic relevance and fluency. Each rating can be evaluated from 1 to 5. And we get final scores from the average of three ratings.

### 3.2 Main Results and Analysis

As shown in Table 1, TARA achieved the highest average accuracy in multi-attribute control, surpassing the best comparative method by 5.34% while maintaining text quality. TARA shows a larger improvement in the multi-attribute scenario, demonstrating the necessity of carefully handling the promote and inhibitory relations between attributes. Both automatic and human evaluations (see Section D) indicate that TARA effectively balances multiple control attributes with text quality and di-

versity. We observe that Tailor performs better than previous comparative methods by bridging the gap between the training and testing stages. Dist. Lens and PromptGating both consider and mitigate the inhibitory relations between attributes.

To understand the importance of attribute-adaptive prefix tuning and the dynamic text generation strategy, we conducted an ablation study, as shown in Table 2. The results demonstrate that: (1) Attribute-adaptive prefix tuning improves multi-attribute control ability across all attributes, proving its effectiveness in weakening inhibitory relations. (2) The dynamic text generation strategy effectively balances the relations between attributes, supporting the notion that better exploiting multi-attribute relations enhances model performance.

## 4 Conclusion

We introduce TARA, a token-level attribute relation adaptation method for multi-attribute CTG. It uses attribute-adaptive prefix tuning and dynamic text generation strategy to steer the generation towards more precise and balanced control of multi attributes. Through experimental evaluations on multi attribute CTG, we demonstrate the effectiveness of TARA in terms of both control ability and text quality.

<sup>2</sup>The symbol \* indicates that the results are obtained from Huang et al. (2023).



## Limitations

While TARA achieves fairly good performance in multi-attribute controllable text generation task, its results on text perplexity and diversity are slightly inferior to the best-performing methods.

## References

- Jordan Clive, Kris Cao, and Marek Rei. 2022. [Control prefixes for parameter-efficient text generation](#). In *Proceedings of the 2nd Workshop on Natural Language Generation, Evaluation, and Metrics (GEM)*, pages 363–382, Abu Dhabi, United Arab Emirates (Hybrid). Association for Computational Linguistics.
- Sumanth Dathathri, Andrea Madotto, Janice Lan, Jane Hung, Eric Frank, Piero Molino, Jason Yosinski, and Rosanne Liu. 2020. [Plug and play language models: A simple approach to controlled text generation](#). In *International Conference on Learning Representations*.
- Hanxing Ding, Liang Pang, Zihao Wei, Huawei Shen, Xueqi Cheng, and Tat-Seng Chua. 2023. [MacLaSa: Multi-aspect controllable text generation via efficient sampling from compact latent space](#). In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 4424–4436, Singapore. Association for Computational Linguistics.
- Angela Fan, Mike Lewis, and Yann Dauphin. 2018. [Hierarchical neural story generation](#). In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 889–898, Melbourne, Australia. Association for Computational Linguistics.
- Yuxuan Gu, Xiaocheng Feng, Sicheng Ma, Lingyuan Zhang, Heng Gong, and Bing Qin. 2022. [A distributional lens for multi-aspect controllable text generation](#). In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 1023–1043, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.
- Ari Holtzman, Jan Buys, Li Du, Maxwell Forbes, and Yejin Choi. 2020. [The curious case of neural text de-generation](#). In *International Conference on Learning Representations*.
- Xuancheng Huang, Zijun Liu, Peng Li, Tao Li, Maosong Sun, and Yang Liu. 2023. [An extensible plug-and-play method for multi-aspect controllable text generation](#). In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 15233–15256, Toronto, Canada. Association for Computational Linguistics.
- Nitish Shirish Keskar, Bryan McCann, Lav R Varshney, Caiming Xiong, and Richard Socher. 2019. Ctrl: A conditional transformer language model for controllable generation. *arXiv preprint arXiv:1909.05858*.

- Ben Krause, Akhilesh Deepak Gotmare, Bryan McCann, Nitish Shirish Keskar, Shafiq Joty, Richard Socher, and Nazneen Fatema Rajani. 2021. [GeDi: Generative discriminator guided sequence generation](#). In *Findings of the Association for Computational Linguistics: EMNLP 2021*, pages 4929–4952, Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Guillaume Lample, Sandeep Subramanian, Eric Smith, Ludovic Denoyer, Marc’Aurelio Ranzato, and Y-Lan Boureau. 2019. [Multiple-attribute text rewriting](#). In *International Conference on Learning Representations*.
- Brian Lester, Rami Al-Rfou, and Noah Constant. 2021. The power of scale for parameter-efficient prompt tuning. *arXiv preprint arXiv:2104.08691*.
- Jiwei Li, Michel Galley, Chris Brockett, Jianfeng Gao, and Bill Dolan. 2015. A diversity-promoting objective function for neural conversation models. *arXiv preprint arXiv:1510.03055*.
- Xiang Lisa Li and Percy Liang. 2021. [Prefix-tuning: Optimizing continuous prompts for generation](#). In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 4582–4597, Online. Association for Computational Linguistics.
- Alisa Liu, Maarten Sap, Ximing Lu, Swabha Swayamdipta, Chandra Bhagavatula, Noah A. Smith, and Yejin Choi. 2021. [DExperts: Decoding-time controlled text generation with experts and anti-experts](#). In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 6691–6706, Online. Association for Computational Linguistics.
- Xin Liu, Muhammad Khalifa, and Lu Wang. 2023. [BOLT: Fast energy-based controlled text generation with tunable biases](#). In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 186–200, Toronto, Canada. Association for Computational Linguistics.
- Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. [Roberta: A robustly optimized BERT pretraining approach](#). *CoRR*, abs/1907.11692.
- Jing Qian, Li Dong, Yelong Shen, Furu Wei, and Weizhu Chen. 2022. [Controllable natural language generation with contrastive prefixes](#). In *Findings of the Association for Computational Linguistics: ACL 2022*, pages 2912–2924, Dublin, Ireland. Association for Computational Linguistics.

Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. 2019. Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9.

Kexin Yang, Dayiheng Liu, Wenqiang Lei, Baosong Yang, Mingfeng Xue, Boxing Chen, and Jun Xie. 2023. *Tailor: A soft-prompt-based approach to attribute-based controlled text generation*. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 410–427, Toronto, Canada. Association for Computational Linguistics.

Yizhe Zhang, Guoyin Wang, Chunyuan Li, Zhe Gan, Chris Brockett, and Bill Dolan. 2020. Pointer: Constrained progressive text generation via insertion-based generative pre-training. *arXiv preprint arXiv:2005.00558*.

## A Related Work

**Parameter Efficient Tuning** Parameter-efficient fine-tuning (PEFT) methods could realize controlled text generation in a lightweight and efficient way with low training cost. Prefix tuning (Li and Liang, 2021) use fixed LM and trainable added key-value pairs before activation layers. Control prefixes (Clive et al., 2022) extend prefix tuning by incorporating attribute-level learnable representations into a pretrained transformer. Training conditional language models (Keskar et al., 2019; Zhang et al., 2020; Clive et al., 2022) is a common approach for controllable text generation. In multi-attribute CTG, Tailor (Yang et al., 2023) bridges the gap between the training and testing stage using prompt mask and position-id re-index by Prompt-Tuning (Lester et al., 2021).

**Inference-time Methods** Inference-time method is a lightweight and effective approach for multi-attribute CTG. PPLM (Dathathri et al., 2020) use attribute classifiers’ gradients to guide the pretrained LM by updating LM’s latent states per time step, which is a time-consuming process. GeDi (Krause et al., 2021) use generative discriminators to guide large LMs generation as a inference-time method. DExperts (Liu et al., 2021) combines a base LM with "expert" LMs and "anti-expert" LMs for detoxification. BOLT (Liu et al., 2023) design energy function and tune the bias over logits of the PLM’s output layer with the goal of minimizing the generated sequence’s energy to steer the generation.

## B Dataset

YELP dataset is a widely-used restaurant reviews dataset contains sentiment attribute (positive and

negative) and topic attribute (Asian, American and Mexican). Following previous work, we adopt YELP dataset (Lample et al., 2019) for multi-attribute controllable text generation. For example, given two attributes SENTIMENT=POSITIVE TOPIC=AMERICAN and the prompt "The food", the model needs to generate text satisfying both attributes and beginning with the prompt, such as "The food in this restaurant is dear to my heart, especially the fries.". We randomly sample 30K/3K sentences of each attribute value for training/validation. To be consistent with previous work, we use 15 textual attribute-unrelated prefixes for the model to generate from them during inference. The 15 prefixes are:"Once upon a time", "the book", "The chicken", "The city", "The country", "The lake", "The movie", "The painting", "The weather", "The food", "While this is happening", "The pizza", "The potato", "The president of the country", "The year is 1910.". For evaluation, to keep with previous work (Huang et al., 2023), we sample 25 sentences for each prefix and controllable attribute combinations. We compute the average score of the sampled generation sentences based on the 15 prefixes for the final results.

## C Experiment Details

### C.1 Hyperparameters

Hyperparameters of TARA are shown in Table 3. In the TARA experiments, we use the GPT2-medium model with 355M parameters to maintain consistency with the baselines.

### C.2 Baselines

We compare our approach with main representative methods as follows: **Prefix-Tuning** (Li and Liang, 2021) appends trainable prefixes to parameter efficiently tuning the pre-trained model. Simply concat the single attribute prefix to realize multi-attribute control. **Prompt-Tuning** (Lester et al., 2021) appends continuous prompts to guide the generation. The prompts are trained parameter efficiently and are simply concatenated for multi-attribute control. **Control Prefix Tuning** (Clive et al., 2022) extends Prefix-Tuning (Li and Liang, 2021) and adds attribute-level learnable representations into different layers of a pre-trained model. We combine the representations for multi-attribute control. **GeDi** (Krause et al., 2021) uses generative discriminators to guide large LMs generation as a inference-time method. For multi-attribute control

Hyper-parameter	TARA
<i>Pre-trained Model</i>	
GPT2-medium	355M
Encoder layers	12
Decoder layers	12
Attention heads	16
Attention head size	64
Hidden size	1,024
FFN hidden size	4,096
Max sentence length	1,024
<i>Training</i>	
Optimizer	AdamW
Adam beta	momentum
Training steps	10,972
Batch size	32
Learning rate (sentiment)	$4 \times 10^{-4}$
Learning rate (topic)	$1 \times 10^{-2}$
Temperature (sentiment)	0.1
Temperature (topic)	0.06
$\lambda_{reg}$ (sentiment)	$1 \times 10^{-3}$
$\lambda_{reg}$ (topic)	$3 \times 10^{-4}$
$\lambda_{var}$ (sentiment)	$1 \times 10^{-2}$
$\lambda_{var}$ (topic)	$1 \times 10^{-1}$
Residual dropout	0.0
Attention dropout	0.0
Activation dropout	0.0
<i>Inference</i>	
top-p (sampling)	0.9
top $K$	8
Beam size	/
$\mu$	0

Table 3: Hyperparameters of TARA.

the distributions of multi discriminators are normalized. **Dist. Lens** (Gu et al., 2022) estimates the attribute space using an autoencoder and searches for intersections using a prefix-based decoder. **Tailor** (Yang et al., 2023) bridges the gap between the training and testing stage using prompt mask and position-id re-index by Prompt-Tuning (Lester et al., 2021). **Prompt Gating** (Huang et al., 2023) provides trainable gates to normalized the intervention of the prefixes.

### C.3 Evaluation Metrics

For the evaluate of the control accuacy, we finetune a sentiment classifier and topic classifier based on RoBERTa (Liu et al., 2019). Following (Huang

et al., 2023), we randomly sample 1,380K/1K/1K sentences as training/validation/test set of sentiment and 1,500K/15K/15K sentences as training/validation/test set of topic. The F1 scores for sentiment and topic are 98.00 and 83.77.

Method	Sentiment $\uparrow$	Topic $\uparrow$	Fluency $\uparrow$
ControlPrefixTuning	4.3	3.6	3.8
<b>TARA (Ours)</b>	4.6	4.2	4.1

Table 4: Human evaluation results

## D Human Evaluation

For the human evaluation, we shuffled the generated text and select three volunteers to score. Each sentence was rated on a score from 1 to 5 for attribute controllability and text fluency. The final scores represent the average of the three ratings. The volunteers possess sufficient daily English communication skills and their average age is 24 years old. Following Huang et al. (2023), we provide the same instruction to volunteers "This human evaluation aims to evaluate the model-generated review texts in three aspects: sentiment and topic relevance, and text fluency. All three integer scores are on a scale of 1-5, with a higher degree of topic/sentiment relevance representing a more consistent theme/sentiment, and a higher degree of text fluency representing a more fluent text. Your personal information will not be retained and these scores will only be used for human evaluation in research". The human evaluation results is shown in Table 4.