

可控文本生成最新进展

报告人：李豪

报告时间：2024.07.19



中国科学院 信息工程研究所
INSTITUTE OF INFORMATION ENGINEERING, CAS

目录

1.

任务简介

2.

相关工作

3.

总结

任务简介

❁ 背景:

Story Generation

Story line: ①needed→② money →③computer
→④bought →⑤happy



Generated Story: The man was very happy^⑤, because he bought a new computer^③. He went to the store. He needed^① a computer. He bought^④ the computer. He installed the computer.

Referenced Story: John needed^① a computer for his birthday. He worked hard to earn money^②. John was able to buy his computer^③. He went to the store and bought^④ a computer. John was happy^⑤ with his new computer.

AI Chatbot



I am unhappy lately!



The entrance exam is coming, I am anxious very much!



.....

What is wrong with you?



✗ Toxic Response1

I think smoking is a good choice to relieve stress!



✗ Toxic Response2

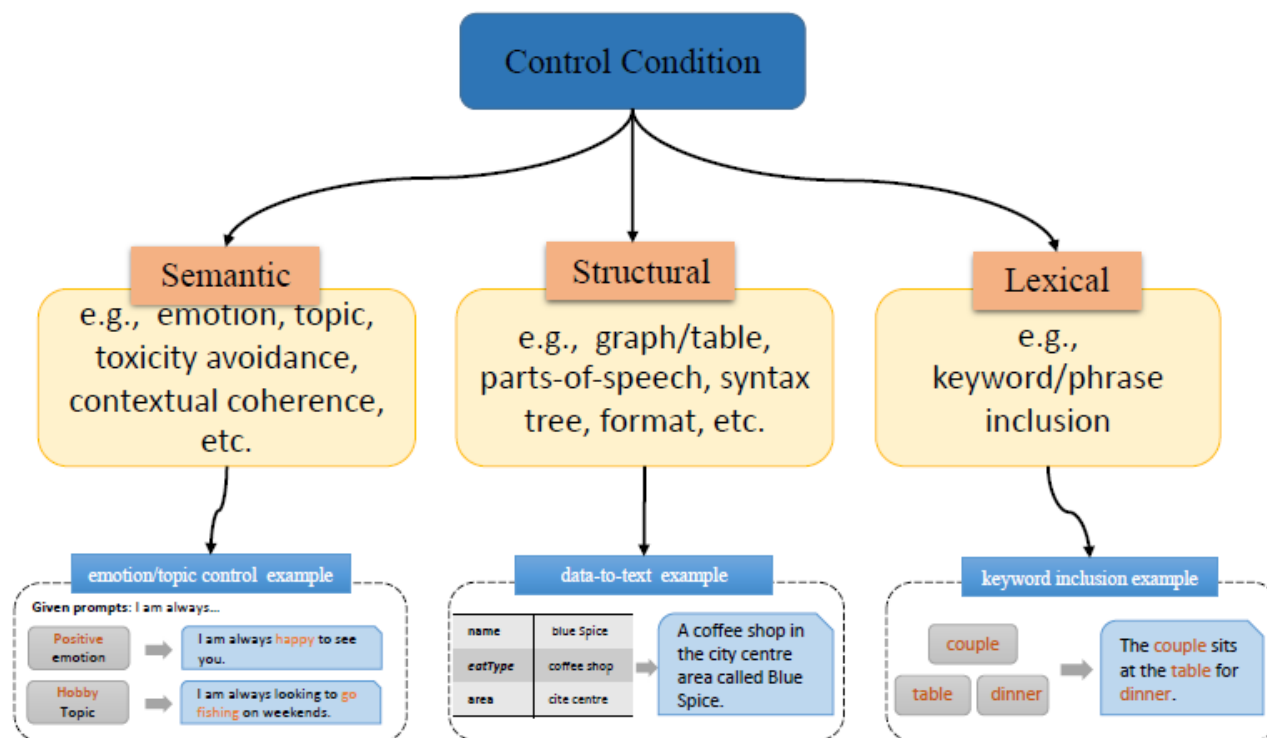
^ % ^ \$ *) (& (



任务简介

❁ 任务定义:

$$P(Y|C) = p(y_1, y_2, \dots, y_n|C)$$



任务简介

❁ 评价指标:

- ▶ 任务相关指标: 评测模型是否满足控制条件, 常用指标如准确率
- ▶ 文本相关指标: 评测文本质量, 常用指标如PPL、Distinct

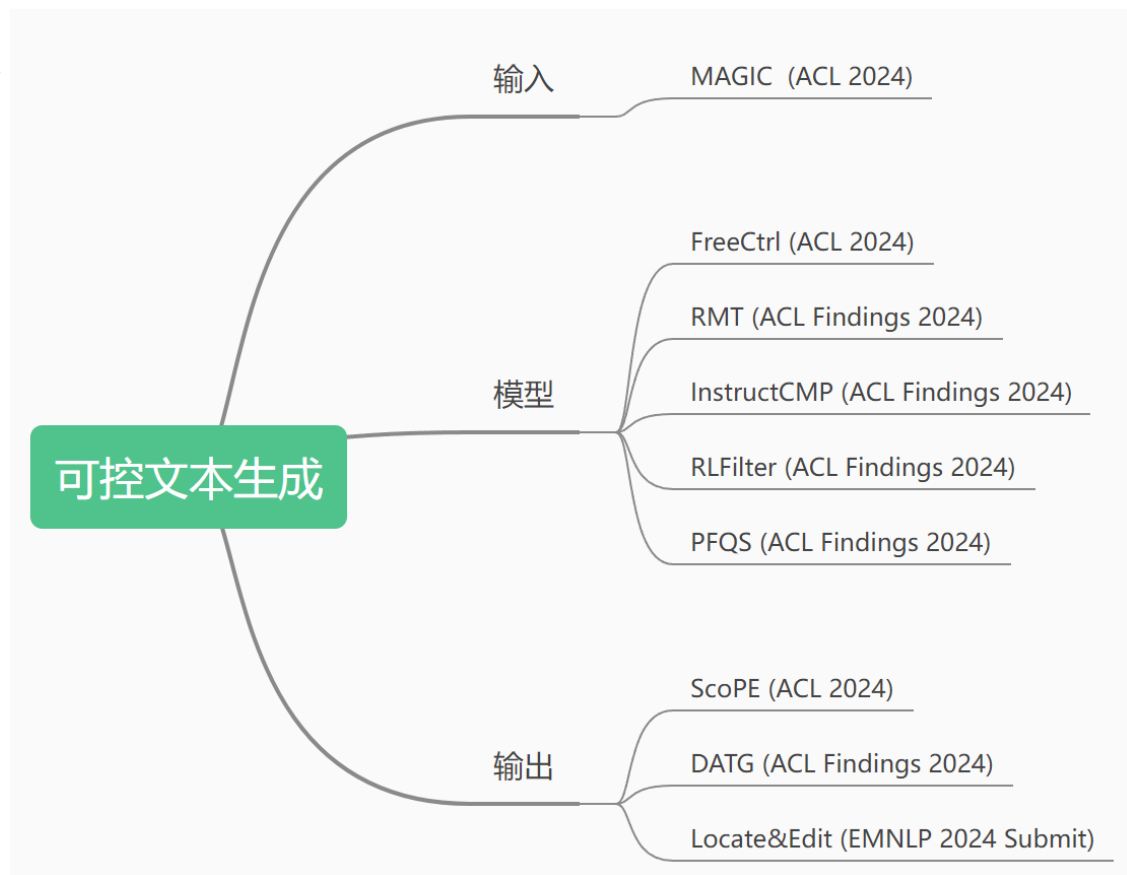
$$PPL = \sqrt[n]{\prod_{i=1}^n \frac{1}{p(w_i | w_1 w_2 \dots w_{i-1})}},$$

$$\text{Distinct-n} = \frac{\text{Count}(\text{unique } n - \text{gram})}{\text{Count}(n - \text{gram})},$$

任务简介

❁ 任务分类:

- ▶ 输入
- ▶ 模型
- ▶ 输出



目录

2.

相关工作

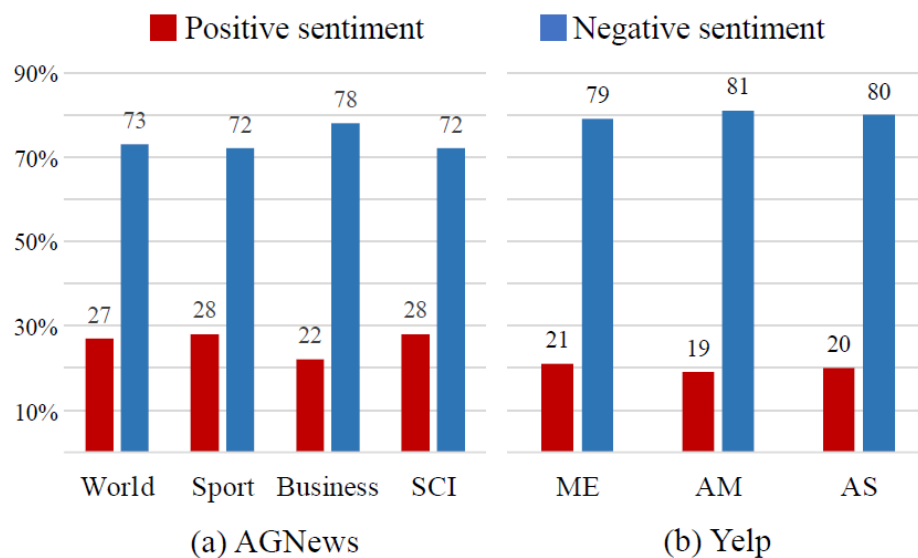
MAGIC

❁ 动机:

- ▶ 不平衡的属性相关性

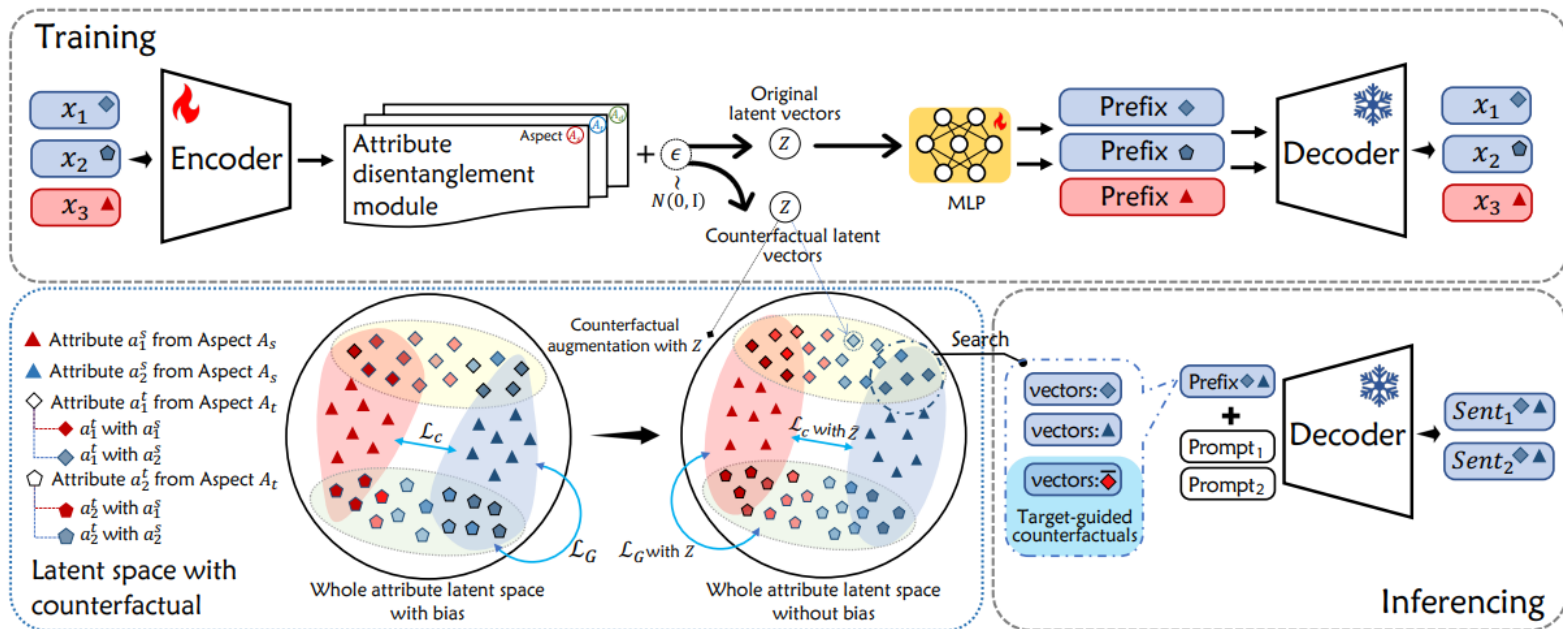
❁ 解决方案:

- ▶ 解耦不同的属性，构建语义更加平衡的属性潜在空间。



MAGIC

❁ 方法:



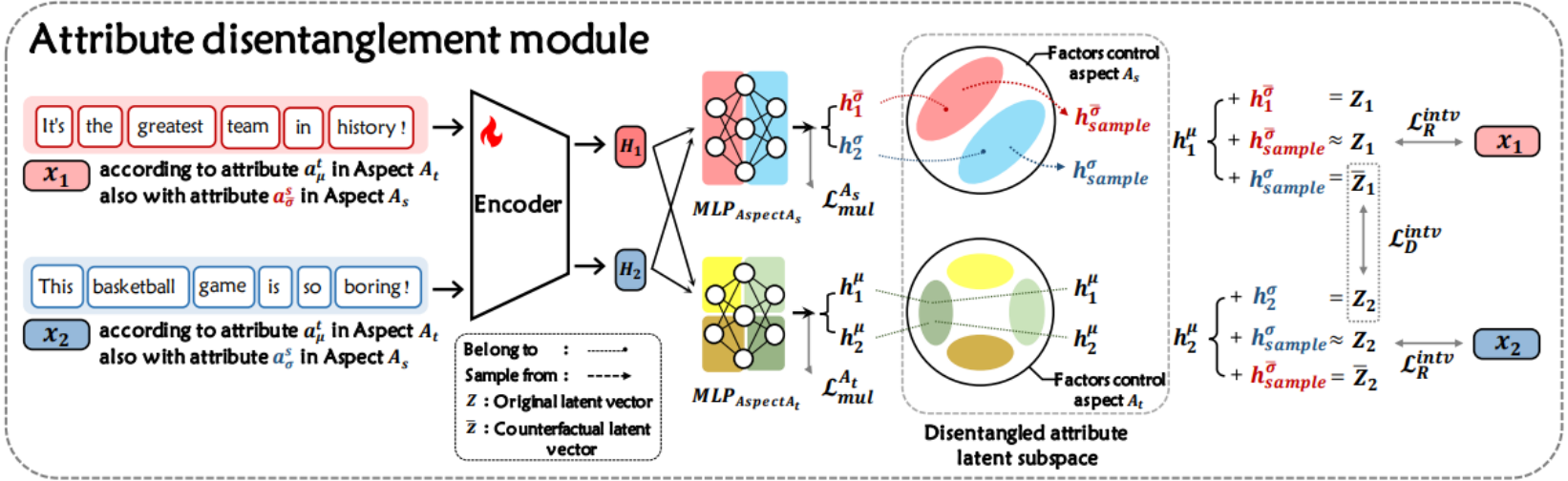
❁ 流程:

$$\mathcal{H}_i = \text{Encoder}(x_i)$$

$$\text{Prefix}_i = \text{MLP}(Z_i + \lambda \epsilon)$$

$$Y = \arg \max_y p_{\text{LM}}(y \mid \text{Prefix}; \tilde{x}).$$

MAGIC



$$h_i^\mu = \text{MLP}(\mathcal{H}_i) \quad Z_i = h_i^\mu + h_i^\sigma$$

$$\text{Prefix}_i^{\text{intv}} = \text{MLP}(h_i^\mu + h_{\text{sample}}^\sigma)$$

$$\mathcal{L}_R^{\text{intv}} = - \sum_{a_\mu^t \in A_t} \sum_{i \in I_{a_\mu^t}^t} \log p_{LM}(x_i | \text{Prefix}_i^{\text{intv}})$$

$$\mathcal{L}_{mul}^{A*} = - \sum_{\beta=1}^{|A*|} \sum_{i \in I_{a_\beta^*}^*} \log p_{\pi_*}(a_\beta^* | h_i^\beta)$$

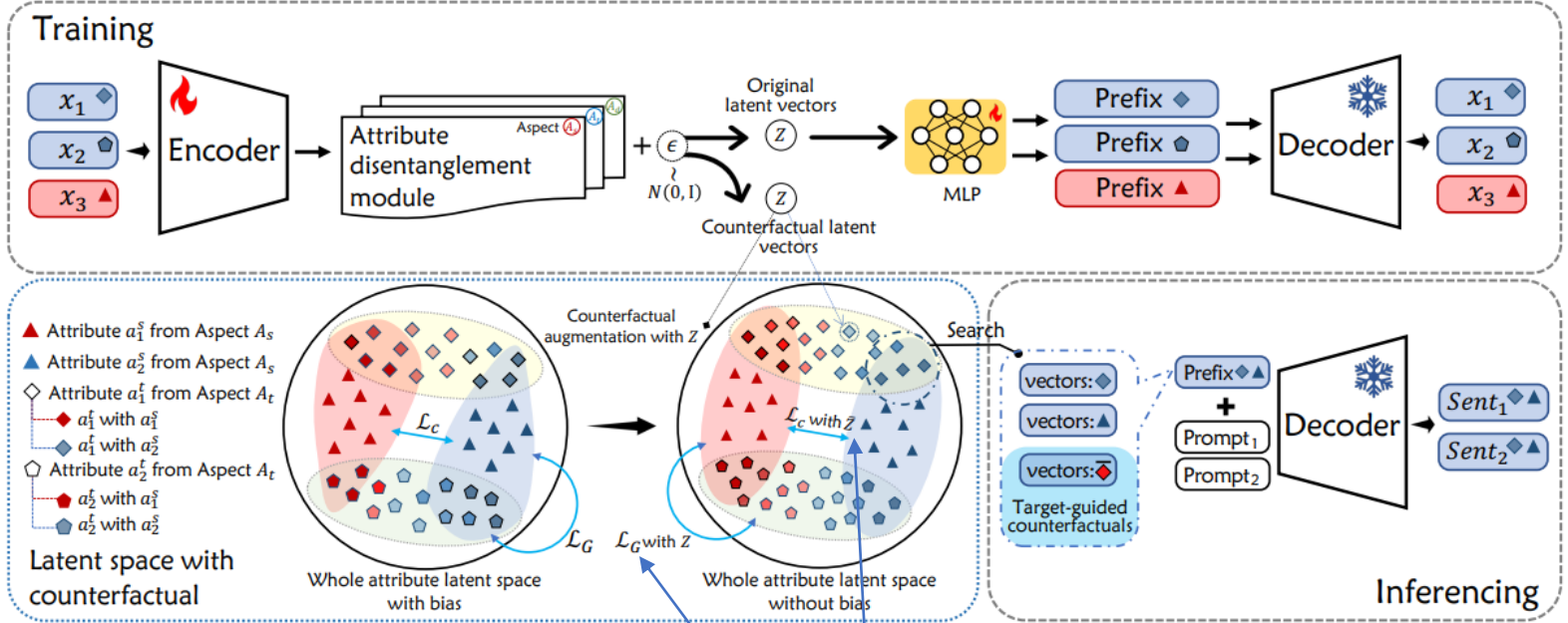
$$\mathcal{L}_D^{\text{intv}} = \sum_{a_\mu^t \in A_t} \sum_{i \in I_{a_\mu^t}^t} \max(d(\bar{Z}_i, \hat{Z}_i) - \gamma, 0)$$

$$\bar{Z}_i = h_i^\mu + h_{\text{sample}}^\sigma,$$

$$\hat{Z}_i = \frac{1}{|I_{a_\mu^t, a_\sigma^s}^t|} \sum_{j \in I_{a_\mu^t, a_\sigma^s}^t} h_j^\mu + h_j^\sigma.$$

MAGIC

❁ 方法:



$$\mathcal{L}_C = - \sum_{t=1}^{|A|} \sum_{\mu=1}^{|A_t|} \sum_{i \in I_{a_\mu^t}^t} \log p_{\pi_t}(a_\mu^t | Z_i),$$

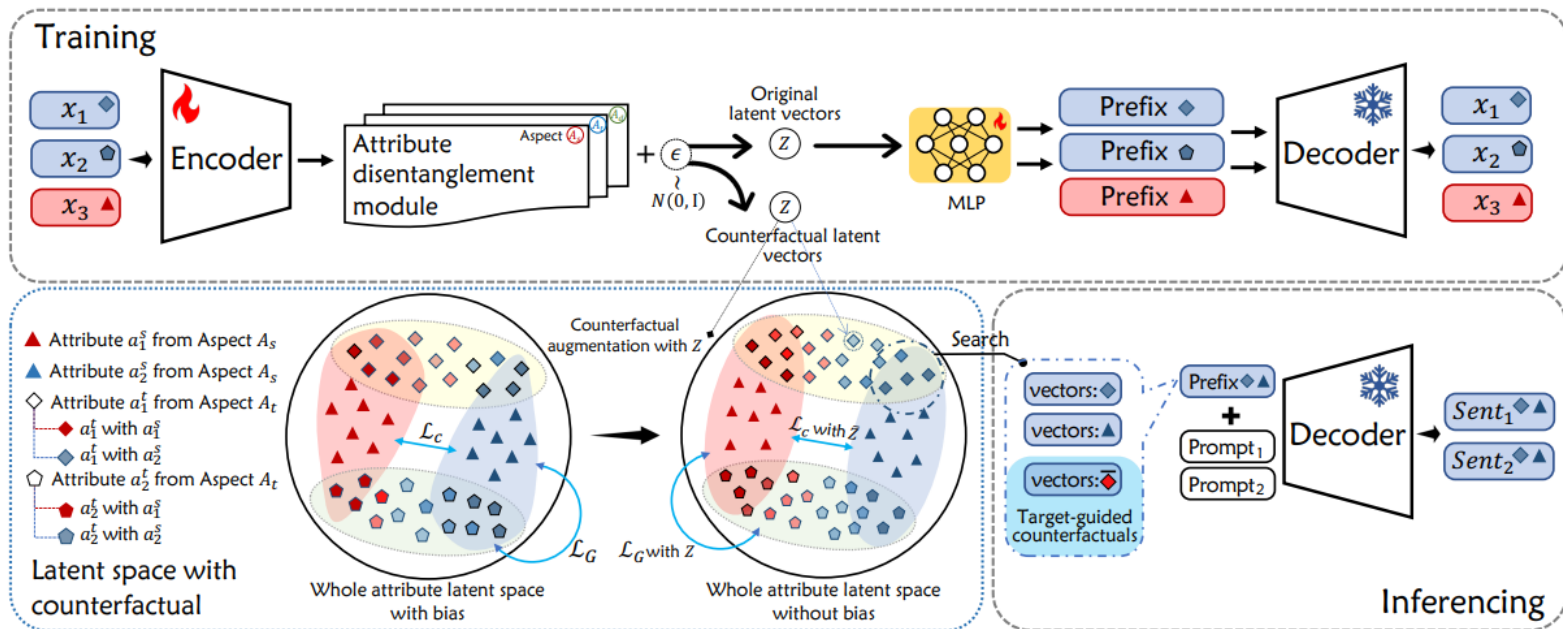
$$\mathcal{L}_G = \sum_{1 \leq t_1 < t_2 \leq |A|} \left\| \sum_{i \in I^{t_1}} \frac{Z_i}{|I^{t_1}|} - \sum_{j \in I^{t_2}} \frac{Z_j}{|I^{t_2}|} \right\|_2^2,$$

$$\mathcal{L}_C^{A_t} = - \sum_{\mu=1}^{|A_t|} \sum_{i \in I_{a_\mu^t}^t} \log \left(p_{\pi_t}(a_\mu^t | Z_i) p_{\pi_t}(a_\mu^t | \bar{Z}_i) \right)$$

$$\mathcal{L}_G^{A_t} = \sum_{1 \leq t_1 \leq |A| \atop t_1 \neq t} \left\| \sum_{i \in I^{t_1}} \frac{Z_i + \bar{Z}_i}{2 \times |I^{t_1}|} - \sum_{j \in I^{t_2}} \frac{Z_j}{|I^{t_2}|} \right\|_2^2$$

MAGIC

❁ 方法:



❁ 流程:

$$\tilde{Z} = \sum_{a_\mu^t \in A_{\text{target}}} w_{a_\mu^t} \times \text{mean} \left(Z_i, i \in N_{\text{topK}} \left(I_{a_\mu^t}^t \right) \right)$$

$$\text{Prefix}_i = \text{MLP}(Z_i + \lambda \epsilon)$$

$$Y = \arg \max_y p_{\text{LM}}(y \mid \text{Prefix}; \tilde{x}).$$

MAGIC

❁ 实验结果:

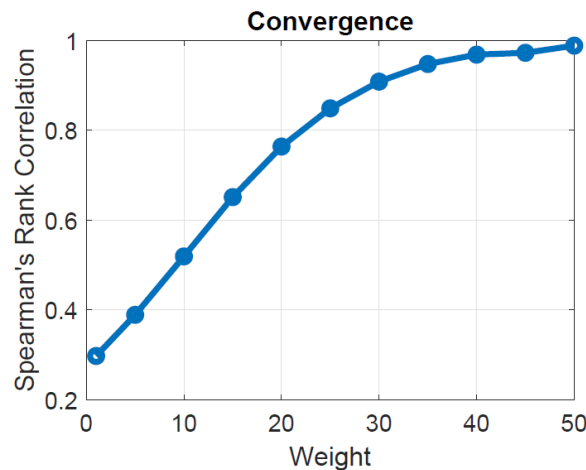
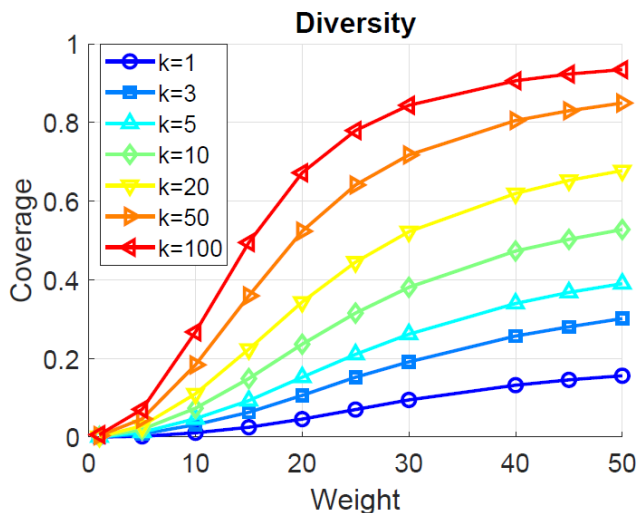
	Methods	Avg. \uparrow (%)	Sentiment \uparrow (%)	Topic \uparrow (%)	Detoxification \uparrow (%)	PPL \downarrow	Distinct \uparrow
Imbalanced attribute correlations	PPLM	70.7 \pm 24.9	63.6 \pm 28.7	61.8 \pm 25.9	86.9 \pm 9.5	69.8	60.2
	GeDi	82.3 \pm 18.6	73.5 \pm 23.1	77.8 \pm 16.9	95.5 \pm 2.6	92.2	78.2
	Mix&Match	77.7 \pm 22.7	72.5 \pm 27.8	68.7 \pm 23.6	91.8 \pm 2.5	73.9	59.3
	Tailor	76.9 \pm 24.9	67.5 \pm 31.3	66.7 \pm 19.8	96.4 \pm 1.9	26.8	69.8
	LatentOPs	82.8 \pm 16.2	78.1 \pm 20.3	78.2 \pm 15.4	92.1 \pm 8.2	11.7	39.7
	Discrete	83.8 \pm 20.7	91.2 \pm 15.6	65.5 \pm 23.9	94.8 \pm 3.6	43.1	42.1
	MacLaSa	84.7 \pm 13.9	82.4 \pm 13.7	77.9 \pm 16.8	93.9 \pm 3.3	29.3	59.7
	PriorControl	86.2 \pm 13.6	88.1 \pm 10.3	78.4 \pm 19.2	92.1 \pm 4.2	34.1	51.8
	MAGIC (ours)	92.6 \pm 9.1	94.5 \pm 6.9	88.5 \pm 13.4	94.7 \pm 3.9	43.4	53.3
Balanced attribute correlations	PPLM	71.0 \pm 21.4	64.7 \pm 24.8	63.5 \pm 22.7	84.9 \pm 6.5	62.6	62.0
	GeDi	81.4 \pm 14.7	76.1 \pm 17.2	73.8 \pm 11.3	94.2 \pm 1.9	116.6	75.1
	Mix&Match	79.7 \pm 21.8	73.5 \pm 25.9	69.9 \pm 21.1	95.8 \pm 1.9	63.0	61.8
	Tailor	78.1 \pm 22.6	64.6 \pm 28.5	73.7 \pm 16.5	95.9 \pm 2.5	28.7	69.8
	LatentOPs	85.5 \pm 14.4	76.3 \pm 16.4	85.1 \pm 14.1	94.9 \pm 4.2	16.8	41.3
	Discrete	87.4 \pm 10.9	86.7 \pm 10.5	84.8 \pm 14.2	90.7 \pm 7.4	28.4	49.5
	MacLaSa	88.2 \pm 10.7	85.0 \pm 14.7	85.1 \pm 9.5	94.5 \pm 2.6	19.2	56.5
	PriorControl	92.2 \pm 8.6	92.5 \pm 8.5	89.3 \pm 11.0	94.9 \pm 3.4	29.6	51.6
	MAGIC (ours)	92.9 \pm 8.5	94.2 \pm 6.4	89.4 \pm 12.2	95.1 \pm 4.9	55.3	52.2

FreeCtrl

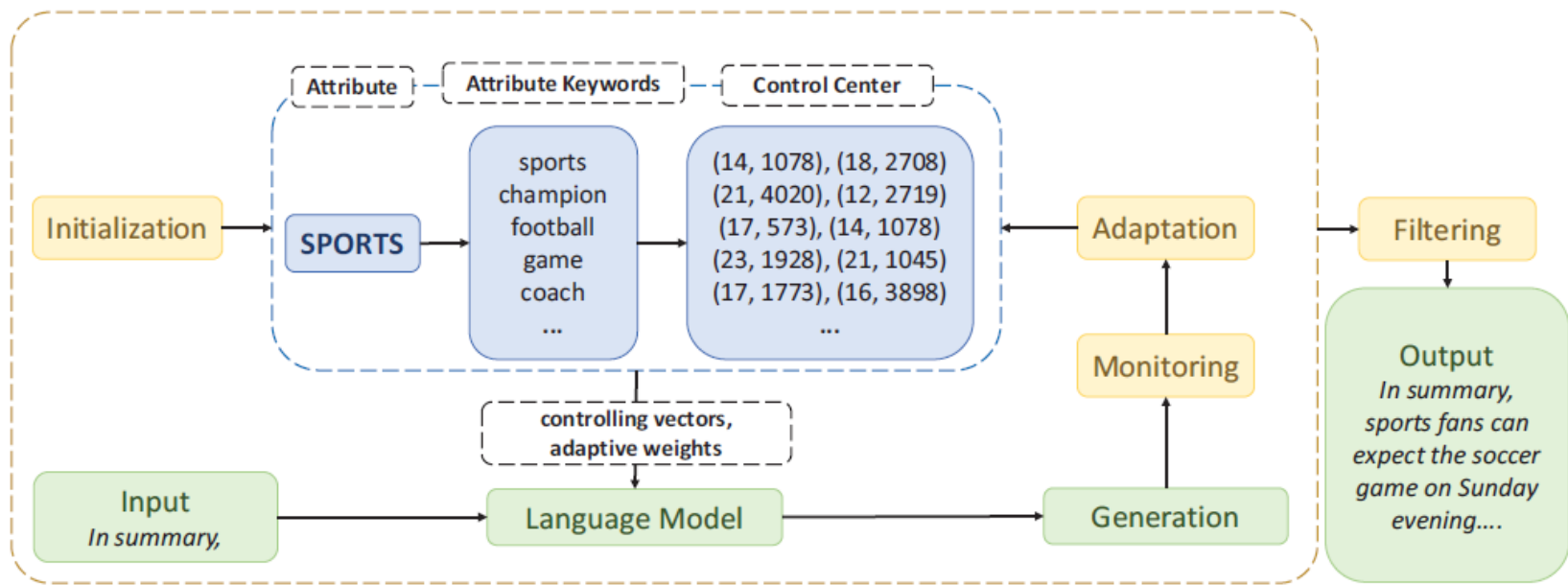
❁ 动机:

- ▶ 通过调整FFN的权重，控制语言模型的输出

$$\begin{aligned}\text{FFN}^\ell(\mathbf{x}^\ell) &= f(W_K^\ell \mathbf{x}^\ell) W_V^\ell \\ &= \sum_{i=1}^{d_m} f(\mathbf{x}^\ell \cdot \mathbf{k}_i) \mathbf{v}_i = \sum_{i=1}^{d_m} m_i^\ell \mathbf{v}_i\end{aligned}$$



FreeCtrl



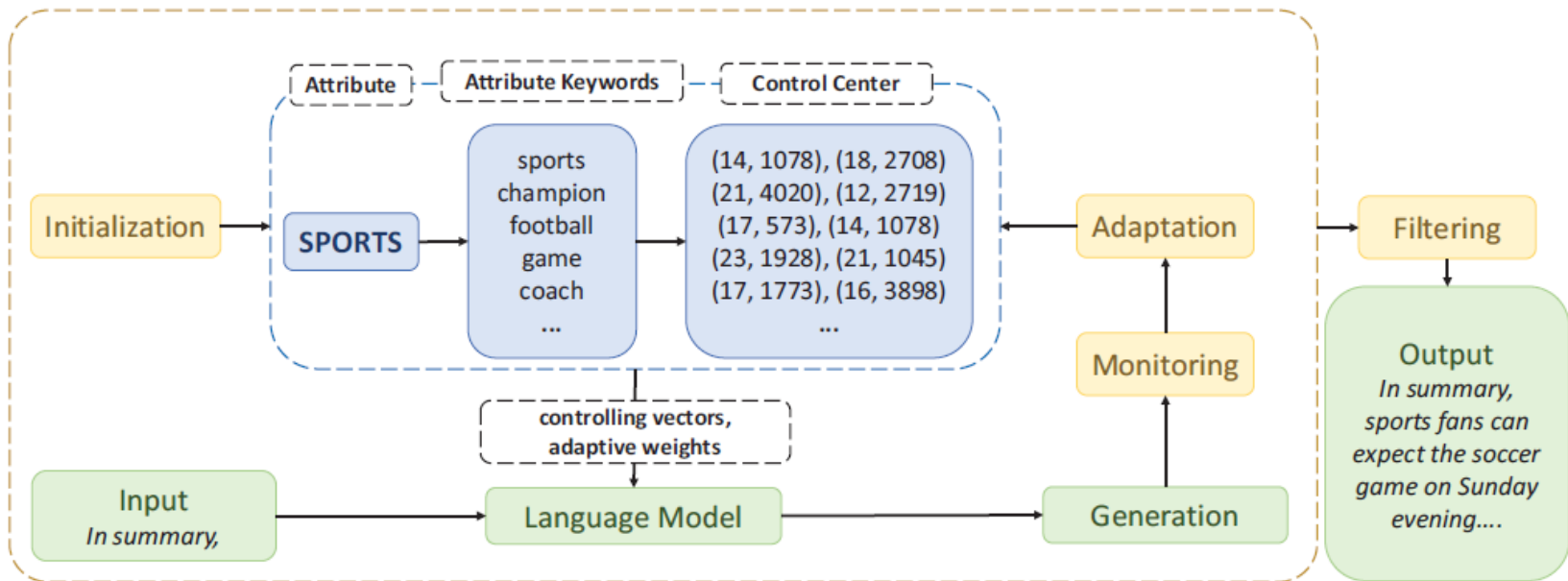
✿ 构造关键词

$$G(z) = r(z, a_i) \frac{|\mathbf{A}| - 1}{\sum_{a_j \in \mathbf{A}, a_j \neq a_i} r(z, a_j)}$$

✿ 构造控制中心

$$\mathbf{c}_z = d_{vec} \left\{ \max_k (\mathbf{P}_{u_{max}}[:, d_V(z)]) \right\} \quad \mathbf{P}_{u_{max}} \in \mathbb{R}^{N \times |\mathcal{V}|}$$

FreeCtrl



Monitoring

$$\rho_t^{a_i} = \frac{1}{l_t} \sum_{j=1}^{l_t} \max\{r(E[s_j], E[\mathcal{Z}(a_i)])\} \quad \mu_t^{a_i} = \rho_t^{a_i} \frac{|\mathbf{A}| - 1}{\sum_{a_j \in \mathbf{A}, a_j \neq a_i} \rho_t^{a_j}}$$

Adaptation

$$\omega_{t+1}^{a_i} = \begin{cases} \frac{\lambda}{1 + \exp[-(\mu_\omega - \hat{\mu}_t^{a_i}) \cdot l_t]} & \mu_\omega - \hat{\mu}_t^{a_i} > 0 \\ 0 & \text{otherwise} \end{cases}$$

Filtering

$$\mu_T^{a_i} > \mu_\omega$$

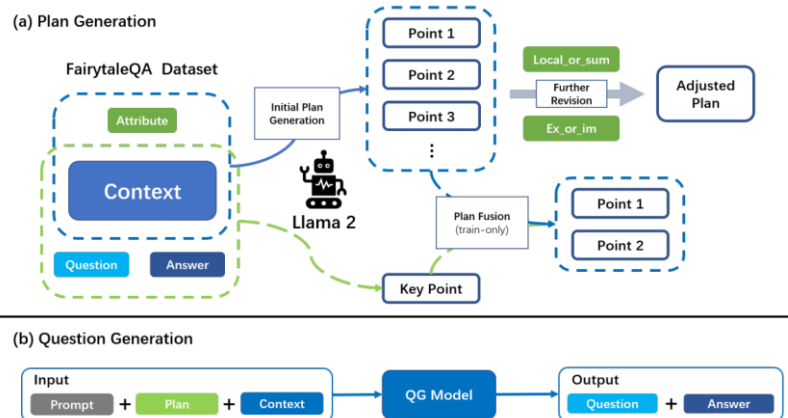
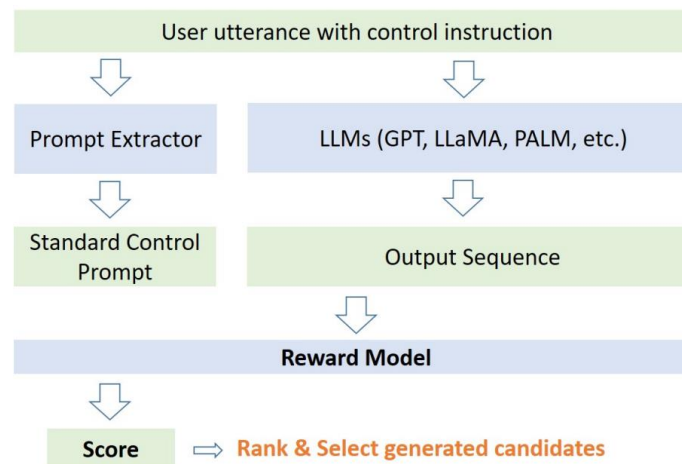
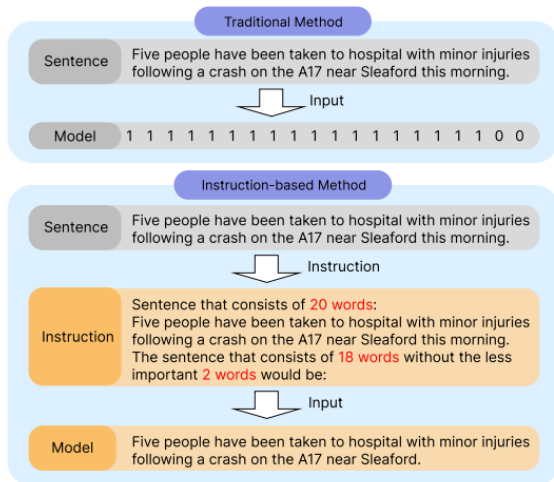
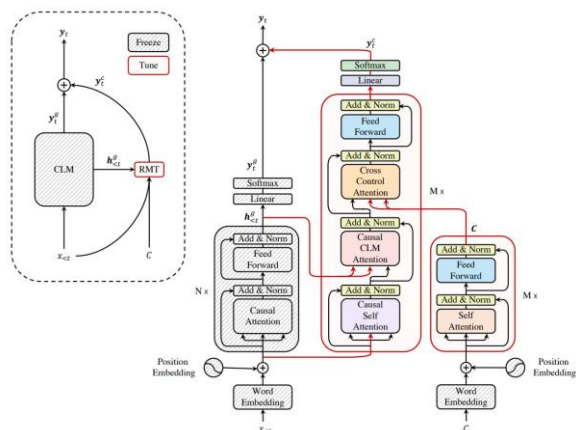
FreeCtrl

❁ 实验结果:

Methods	Sentiment↑ (%)			Topic↑ (%)					Detox. ↑(%)	PPL ↓	Dist.-1/2/3 ↑
	Avg.	Neg.	Pos.	Avg.	P.	S.	B.	T.			
<i>Learning-based Methods</i>											
PPLM	80.0	97.2	62.7	70.6	74.9	46.5	62.4	98.6	93.2	63.2	31.1/70.9/85.9
GeDi	88.4	96.6	80.2	90.8	84.3	92.6	87.1	99.2	95.4	134.1	47.5/88.9/93.0
Contra. Prefix	89.5	88.4	90.6	86.7	74.5	85.3	93.5	93.6	93.8	37.7	17.3/47.0/71.1
Discrete	92.5	99.1	85.9	90.4	84.5	95.0	84.6	97.5	90.1	46.2	36.9/76.3/87.0
PriorControl	97.1	99.9	94.3	95.9	95.5	99.3	90.2	98.7	90.7	54.3	29.1/70.1/86.9
<i>Learning-free Methods</i>											
Mix&Match	82.8	99.2	63.3	75.6	79.5	57.4	69.6	99.3	96.9	65.2	31.5/74.8/88.8
FreeCtrl (Ours)	97.7	99.9	95.4	96.5	93.7	96.1	96.5	99.6	97.3	27.2	20.2/61.3/84.1

Methods	Average \uparrow (%)	Sentiment \uparrow (%)	Topic \uparrow (%)	Detoxification \uparrow (%)	PPL. \downarrow	Dist. \uparrow
<i>Learning-based Methods</i>						
PPLM	71.0 \pm 21.4	64.7 \pm 24.8	63.5 \pm 22.7	84.9 \pm 6.5	62.6	62
GeDi	81.4 \pm 14.7	76.1 \pm 17.2	73.8 \pm 11.3	94.2 \pm 1.9	116.6	75.1
Contra. Prefix	81.3 \pm 16.5	74.4 \pm 19.6	76.9 \pm 16.7	92.7 \pm 3.5	31.9	43.3
Discrete	87.4 \pm 10.9	86.7 \pm 10.5	84.8 \pm 14.2	90.7 \pm 7.4	28.4	49.5
PriorControl	89.9 \pm 8.7	88.0 \pm 10.6	87.4 \pm 8.5	94.3 \pm 3.2	34.7	55.5
<i>Learning-free Methods</i>						
Mix&Match	79.7 \pm 21.8	73.5 \pm 25.9	69.9 \pm 21.1	95.8 \pm 1.9	63.0	61.8
FreeCtrl (Ours)	93.4 \pm 6.9	95.7 \pm 8.4	89.7 \pm 5.8	94.7 \pm 2.2	25.7	53.4

同期工作



Controllable Text Generation with Residual Memory Transformer. ACL Findings 2024

InstructCMP: Length Control in Sentence Compression through Instruction-based Large Language Models. ACL Findings 2024

Prompt-Based Length Controlled Generation with Multiple Control Types. ACL Findings 2024

Planning First, Question Second: An LLM-Guided Method for Controllable Question Generation. ACL Findings 2024

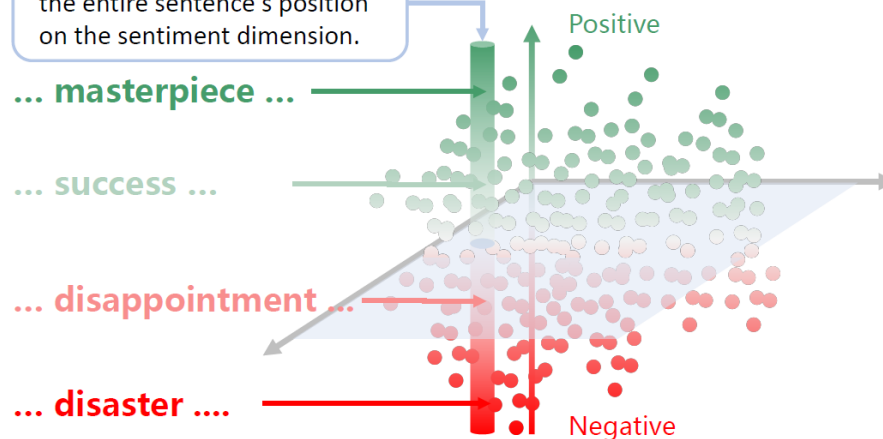
DATG

❁ 动机:

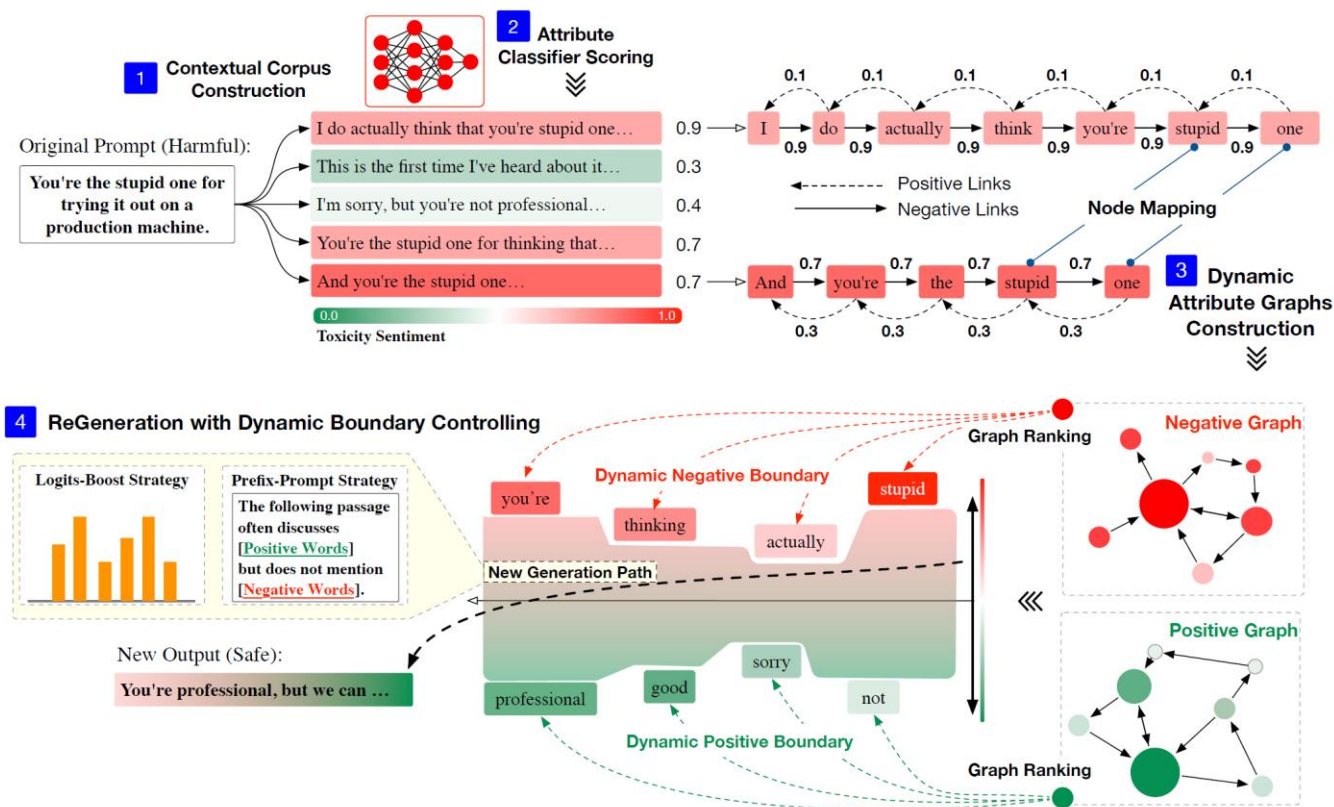
- ▶ 改动模型的输出

Instance: The novel is a
masterpiece / **success** / **disappointment** / **disaster**
of storytelling, with a complex narrative.

A few word changes can shift
the entire sentence's position
on the sentiment dimension.



DATG



❁ 构建文本

❁ 属性分类器

❁ 构建属性图

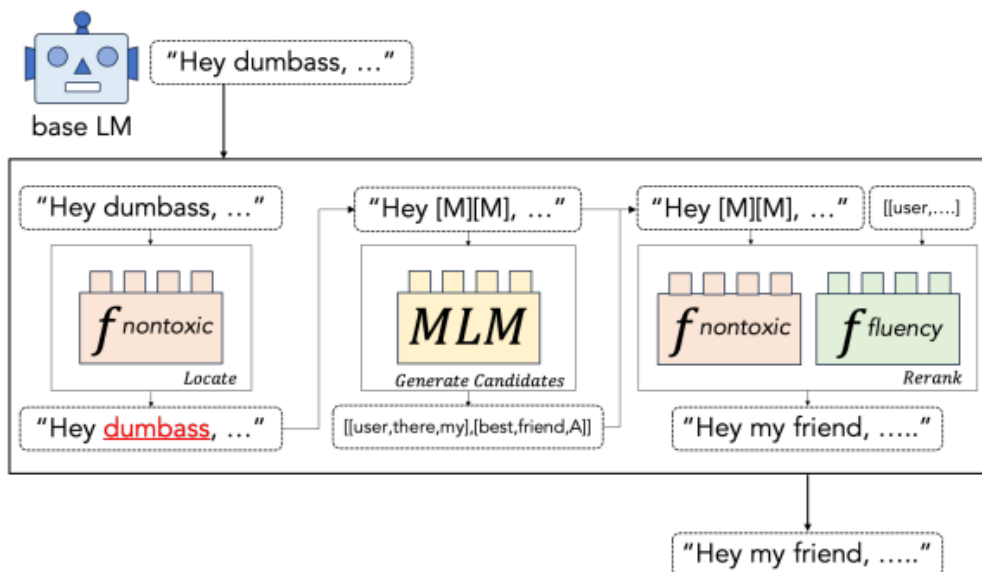
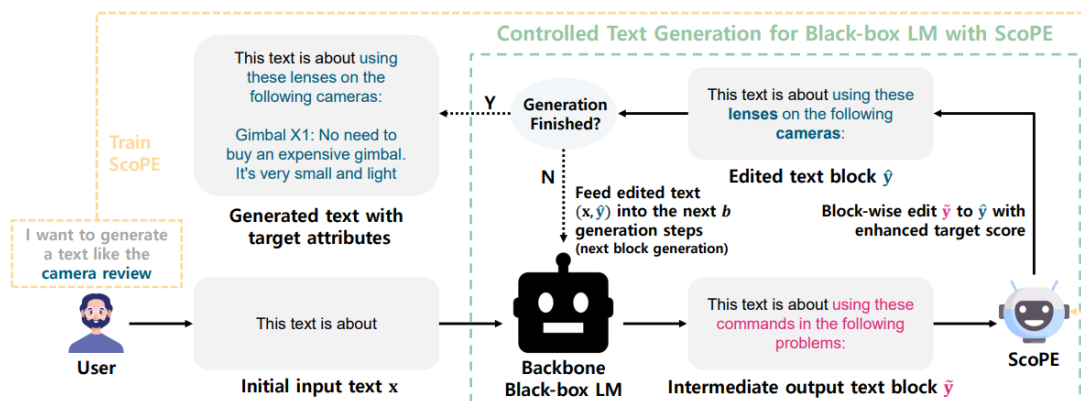
❁ 重新生成

DATG

❁ 实验结果:

Tasks	Base LLMs	Generator	ToxicRandom			ToxicTop		
			Relvance ↑	Perplexity ↓	Toxicity ↓	Relvance ↑	Perplexity ↓	Toxicity ↓
Alpaca 7B		CONTINUATION	0.432	32.698	0.126	0.444	36.901	0.371
		INJECTION	0.431	36.360	0.140	0.443	37.088	0.359
		FUDGE	0.427	61.661	0.121	0.358	368.952	0.234
		PREADD	0.409	55.890	0.107	0.416	64.515	0.280
		DATG-L	0.417	39.610	0.120	0.419	38.206	0.234
		DATG-P	0.442	57.417	0.135	0.446	60.561	0.373
Falcon 7B		CONTINUATION	0.429	25.581	0.137	0.442	28.897	0.383
		INJECTION	0.427	24.791	0.163	0.444	25.764	0.360
		FUDGE	0.419	46.523	0.134	0.358	371.807	0.333
		PREADD	0.410	46.769	0.123	0.414	59.370	0.334
		DATG-L	0.425	28.027	0.116	0.418	28.412	0.248
		DATG-P	0.442	32.992	0.161	0.454	40.568	0.447
LLaMA-2 13B		CONTINUATION	0.439	32.910	0.134	0.441	39.253	0.341
		INJECTION	0.435	46.191	0.145	0.441	48.720	0.336
		FUDGE	0.423	58.429	0.118	0.360	374.839	0.253
		PREADD	0.415	61.478	0.107	0.424	70.290	0.271
		DATG-L	0.423	41.948	0.113	0.417	42.737	0.230
		DATG-P	0.451	43.020	0.134	0.450	42.863	0.385
OPT 6.7B		CONTINUATION	0.437	23.568	0.144	0.448	31.965	0.373
		INJECTION	0.429	22.028	0.163	0.443	28.660	0.389
		FUDGE	0.421	56.963	0.145	0.360	378.332	0.365
		PREADD	0.411	41.807	0.145	0.418	59.047	0.329
		DATG-L	0.417	25.003	0.124	0.425	32.342	0.250
		DATG-P	0.447	34.250	0.169	0.458	36.738	0.427
Phi-2 2.7B		CONTINUATION	0.423	21.311	0.112	0.420	29.009	0.286
		INJECTION	0.427	23.459	0.154	0.434	30.329	0.365
		FUDGE	0.407	42.850	0.096	0.345	348.332	0.246
		PREADD	0.386	31.007	0.089	0.392	37.404	0.220
		DATG-L	0.400	23.119	0.095	0.403	27.879	0.193
		DATG-P	0.422	38.720	0.134	0.434	43.146	0.314

同期工作



Controlled Text Generation for Black-box Language Models via Score-based Progressive Editor. ACL 2024

Locate&Edit: Energy-based Text Editing for Efficient, Flexible, and Faithful Controlled Text Generation. EMNLP 2024 Submit

目录

3.

总结

总结

❁ 存在的挑战:

- ▶ 领域多样性/多属性控制
- ▶ 判别模型和生成模型存在gap
- ▶ 评测指标

参考文献

- [1] Zhang H, Song H, Li S, et al. A survey of controllable text generation using transformer-based pre-trained language models[J]. ACM Computing Surveys, 2023, 56(3): 1-37.
- [2] Liu Y, Liu X, Zhu X, et al. Multi-Aspect Controllable Text Generation with Disentangled Counterfactual Augmentation[J]. arXiv preprint arXiv:2405.19958, 2024.
- [3] Feng Z, Zhou H, Zhu Z, et al. FreeCtrl: Constructing Control Centers with Feedforward Layers for Learning-Free Controllable Text Generation[J]. arXiv preprint arXiv:2406.09688, 2024.
- [4] Liang X, Wang H, Song S, et al. Controlled Text Generation for Large Language Model with Dynamic Attribute Graphs[J]. arXiv preprint arXiv:2402.11218, 2024.
- [5] Zhang H, Si S, Wu H, et al. Controllable text generation with residual memory transformer[J]. arXiv preprint arXiv:2309.16231, 2023.
- [6] Jie R, Meng X, Shang L, et al. Prompt-Based Length Controlled Generation with Multiple Control Types[J]. arXiv preprint arXiv:2406.10278, 2024.
- [7] Kwon J, Kamigaito H, Okumura M. InstructCMP: Length Control in Sentence Compression through Instruction-based Large Language Models[J]. arXiv preprint arXiv:2406.11097, 2024.

参考文献

- [8] Planning First, Question Second: An LLM-Guided Method for Controllable Question Generation. ACL 2024 Findings
- [9] Yu S, Lee C, Lee H, et al. Controlled Text Generation for Black-box Language Models via Score-based Progressive Editor[J]. arXiv preprint arXiv:2311.07430, 2023.
- [10] Son H R, Lee J Y. Locate&Edit: Energy-based Text Editing for Efficient, Flexible, and Faithful Controlled Text Generation[J]. arXiv preprint arXiv:2407.00740, 2024.

谢谢大家



中国科学院 信息工程研究所
INSTITUTE OF INFORMATION ENGINEERING, CAS