

Introduction

Interesting Problem

The story ending generation methods based on **Pre-trained Language Model** and the **MLE training objective** lack the support of causal, temporal, and sentiment knowledge, resulting in **poor consistency** between generated story ending and the story context.

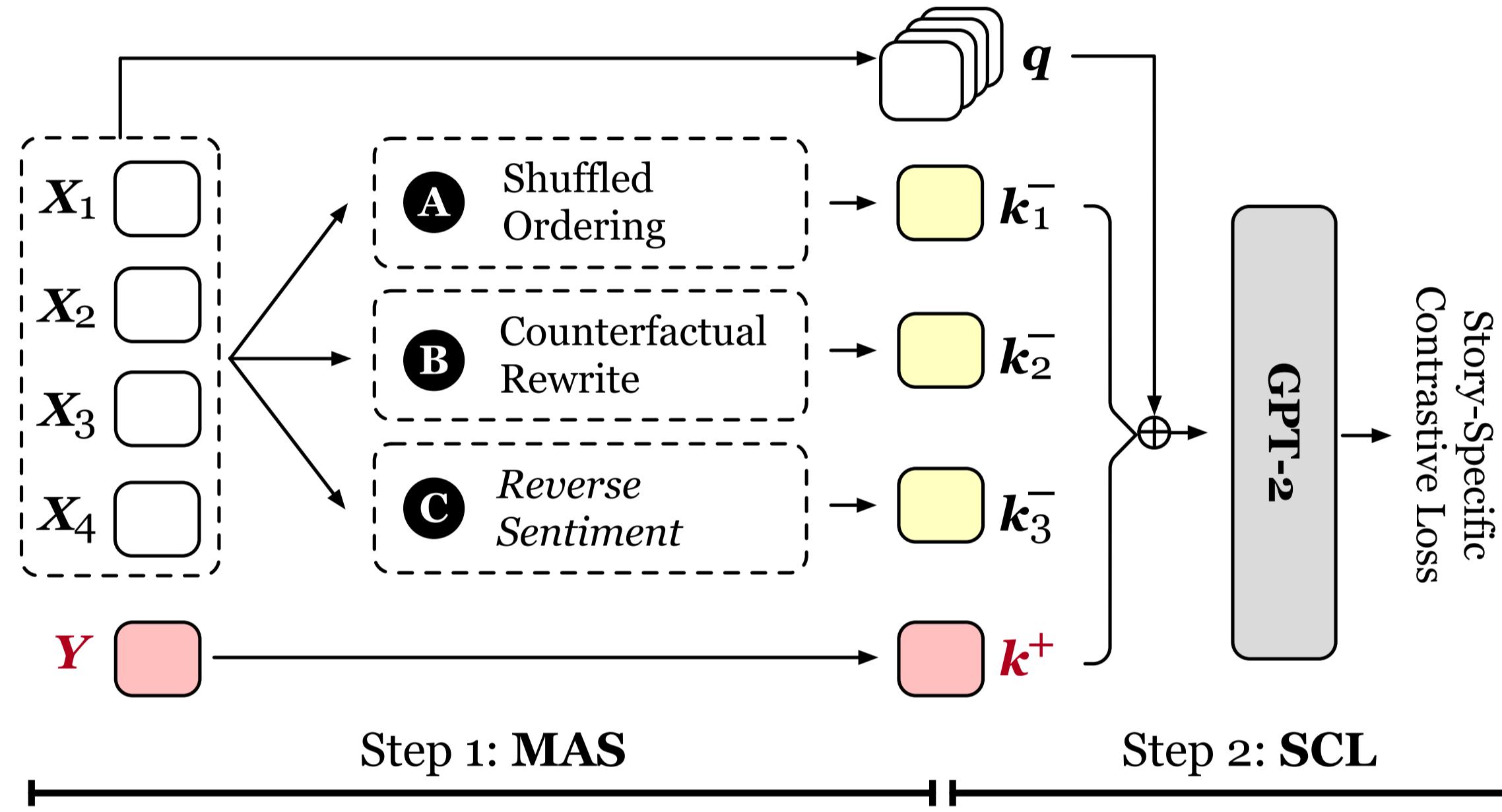
Two crucial challenges in CLSEG:

- How to generate **high-quality wrong endings** inconsistent with story contexts?
- How to design a contrastive training strategy **adapted for story ending generation**?

Contributions:

- The **Multi-Aspect negative sampling** is performed from the perspectives of causality, temporal, and sentiment;
- A **Story-level Contrastive Learning Objective** is designed, which enables the model to compare multi-source knowledge during the training process and improves consistency of the generated story ending;
- Automatic and manual evaluations demonstrate that CLSEG outperforms baselines and can generate story endings with stronger consistency and rationality.

Overview of CLSEG

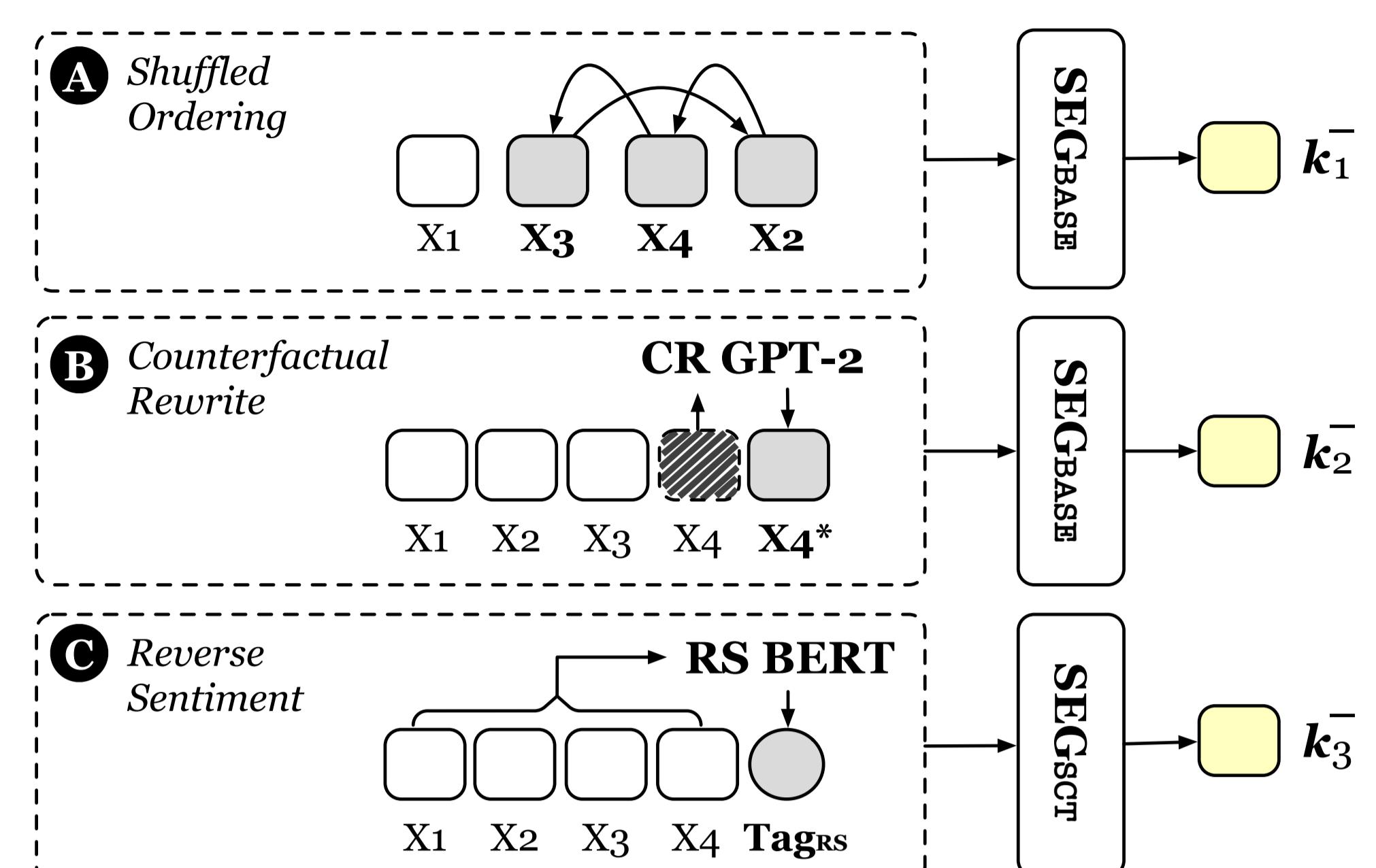


We propose a novel Contrastive Learning framework for Story Ending Generation (CLSEG)†, which has two steps: multi-aspect sampling and story-specific contrastive learning.

Details of CLSEG

Multi-Aspect Sampling

generate high-quality wrong endings



order of events is different -> inconsistent ending

$$\mathcal{Y}_{SO}^- = \text{SEGBASE}(\text{shuffled}(\{\mathcal{X}_1, \mathcal{X}_2, \dots, \mathcal{X}_m\}))$$

what if event A turns B -> inconsistent ending

$$\mathcal{Y}_{CR}^- = \text{SEGBASE}(\{\mathcal{X}_1, \mathcal{X}_2, \dots, \mathcal{X}_m^*\})$$

change the sentiment orientation -> inconsistent ending

$$\mathcal{Y}_{RS}^- = \text{SEGSCT}(\{\mathcal{X}_1, \mathcal{X}_2, \dots, \mathcal{X}_m, \text{Tag}_{RS}\})$$

Notes:

- ◆ **SEGBASE**: GPT-2 fine-tuned on **ROCStories** corpus
- ◆ **SEGSCT**: GPT-2 fine-tuned on selected stories of **SCT**, where wrong ending has reversed sentiment to the story context
- ◆ **CR GPT-2**: GPT-2 fine-tuned on counterfactual story corpus, **TIMETRAVEL**
- ◆ **RS BERT**: BERT_{LARGE} fine-tuned on **SST-2**

Story-level Contrastive Learning Objective

adapted for story ending generation

$$\mathcal{L}_{SCL}^t(p_\theta(\mathcal{Y}|\mathcal{X}), \{\mathcal{Y}_i^-\}_{i=1}^3) = -\log p_\theta(k_t^+ | k_{<t}^+, \mathbf{q}) \\ - \alpha \cdot \frac{1}{N} \sum_i^N \log (1 - p_\theta(k_{i,t}^- | k_{i,<t}^-, \mathbf{q}))$$

where \mathcal{X} denoted as \mathbf{q}

$\{\mathcal{Y}_{SO}^-, \mathcal{Y}_{CR}^-, \mathcal{Y}_{RS}^-\}$ denoted as $\{k_i^-\}_{i=1}^3 = \{k_1^-, k_2^-, k_3^-\}$ for SCL

Experiment & Analysis

Dataset: ROCStories (Train 78,530/Dev 9,816/Test 9,816)

Baselines

- GPT-2(PT)**: vanilla GPT-2.
- GPT-2(FT)**: vanilla GPT-2 fine-tuned on the ROCStories corpus with LM objective.
- GPT-2(GCL)**: This model is pre-trained similar to our model, only the negative sampling is the general noisy ways: [random shuffle/drop/replace tokens](#).

Human Evaluation

Metrics

- Content Quality: whether the generated story ending is fluent and coherent.
- Content Rationality: whether the story endings are consistent with the story context.

Results

Model	Quality \uparrow			Consistent \uparrow	
	Fluency	Coherence	Order	Causal	Sentiment
GPT-2(PT)	1.87	1.58	1.22	1.07	1.19
GPT-2(FT)	2.24	2.13	1.47	1.67	1.68
GPT-2(GCL)	2.08	2.07	1.62	1.73	1.63
CLSEG	2.20	2.43	2.08	2.19	2.02

Table 2: Manual Evaluation in terms of quality and rationality about the generated story endings.

Main Results

Models	BLEU \uparrow	R-1-P \uparrow	R-1-R \uparrow	R-1-F1 \uparrow	R-2-P \uparrow	R-2-R \uparrow	R-2-F1 \uparrow	R-L-P \uparrow	R-L-R \uparrow	R-L-F1 \uparrow	Meteor \uparrow
GPT-2(PT)	1.14	13.99	15.99	13.56	1.34	2.51	1.68	13.17	15.89	12.65	10.36
GPT-2(FT)	2.57	15.15	12.55	13.30	2.12	1.82	1.87	13.87	11.45	12.14	10.48
GPT-2(GCL)	1.71	14.81	14.66	14.21	1.82	2.13	1.89	13.67	12.56	13.12	10.54
CLSEG	1.97	15.68	14.73	14.63	1.87	2.15	1.91	14.44	13.61	13.48	10.73

Table 4: Generated story endings by different models.

Ablation Study

Model	BLEU \uparrow	R-1-R \uparrow	R-2-R \uparrow	R-L-R \uparrow	Meteor \uparrow
CLSEG	1.97	14.73	2.15	13.61	10.73
only SO	1.94	13.77	1.96	12.78	10.16
only CR	2.13	13.71	1.96	12.66	10.20
only RS	2.33	14.24	2.03	13.02	10.64

Table 3: Ablation study of CLSEG.

Error Analysis

Error Type	Story Context + Ending
Repetition	Morgan enjoyed long walks on the beach. ... Morgan decided to propose to her boyfriend. <i>The walk to the beach and the propose.</i>
Conflicting	Frank ... Since Frank was already a bit drunk, he could not drive. <i>Frank had to drive to his date.</i>
Ambiguous	Sunny enjoyed going to the beach. As she stepped out of her car, ... Sunny got back into her car and heading towards the mall. <i>She was going to the park.</i>

Table 5: Typical errors generated by CLSEG. *Italic* words denote the error generated story endings.

Acknowledgments

We thank all anonymous reviewers for their constructive comments. This work is supported by the National Natural Science Foundation of China (No. 62006222 and No. U21B2009).