

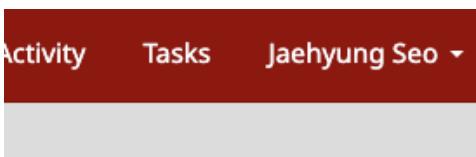
OpenRe

← Back to

CoM  
Editi



# E: An Unlearning-based Approach to Conflict-free M nc



Model



- Counterfa

---

**GPT-J**

---

Full-FT

F-Learni

OURS-m

OURS-pn

---

**LLaMA-3**

---

Full-FT

F-Learni

OURS-m

OURS-pr

- ZsRE Data

---

**GPT-J**

---

Full-FT

F-Learni

OURS-m

OURS-pr

---

**LLaMA-3**

---

Full-FT

F-Learni

OURS-m

OURS-pr

Once again, we

**Format Chec****Copyright PD****Handbook**

## Knowledge Editing, Unlearning, Locate-then-Edit Line Learning for NLP

Language models (LLMs) often retain outdated or incorrect information from pre-training, which undermines their reliability. While model editing methods can mitigate this issue through full re-training, they frequently suffer from knowledge conflicts, where outdated information interferes with new knowledge. In this work, we propose CoME, a framework that enhances the accuracy of knowledge updates in LLMs by selectively removing outdated knowledge. CoME leverages unlearning to mitigate knowledge conflicts while maintaining relevant linguistic features. Through experiments on GPT-J and LLaMA-3 using Counterfactual and ZsRE datasets, we show that CoME significantly improves the accuracy and model reliability when applied to existing editing methods. Our results highlight that the targeted removal of outdated knowledge is critical for maintaining the model's generative performance.

[/forum?id=XCks7AD5HR](#)

### Metareview:

thank you for reviewing our paper and for acknowledging the contribution of our method in alleviating knowledge conflict in knowledge editing. Your thoughtful review and the detailed feedback provided in this round will fully address the concerns raised by the reviewers.

We would like to share the FT experiment results that we could not provide due to the limited rebuttal period. The results highlight that our approach outperforms F-learning in terms of both efficacy and generality. Specifically, F-learning—a knowledge editing approach based on FT—attempts to mitigate conflicts by erasing previously learned knowledge. While this approach can effectively alleviate conflicts, it is inherently limited by the constraints of the FT framework, leading to suboptimal performance. In contrast, our method can effectively mitigate knowledge conflicts while utilizing resources more efficiently. These findings will be incorporated into the final version of the manuscript.

### Dataset Results

Score Efficacy Generality Specificity

	Score	Efficacy	Generality	Specificity
CoME	35.6	29.0	28.1	71.4

CoME	38.1	30.5	30.8	<b>73.7</b>
------	------	------	------	-------------

CoME	<b>86.4</b>	99.4	91.1	73.2
------	-------------	------	------	------

CoME	<b>86.4</b>	<b>99.8</b>	<b>95.3</b>	70.3
------	-------------	-------------	-------------	------

3

CoME	28.7	37.5	36.6	62.8
------	------	------	------	------

CoME	32.1	25.0	23.9	<b>84.8</b>
------	------	------	------	-------------

CoME	78.2	<b>95.7</b>	<b>91.3</b>	59.0
------	------	-------------	-------------	------

CoME	<b>82.3</b>	92.4	83.6	73.3
------	-------------	------	------	------

### Dataset Results

Score Efficacy Generality Specificity

CoME	37.4	52.2	49.6	24.5
------	------	------	------	------

CoME	39.8	58.8	55.4	24.8
------	------	------	------	------

CoME	<b>50.3</b>	<b>97.3</b>	<b>93.0</b>	25.9
------	-------------	-------------	-------------	------

CoME	49.0	89.4	83.1	<b>26.3</b>
------	------	------	------	-------------

3

CoME	54.1	61.5	60.1	44.2
------	------	------	------	------

CoME	55.5	64.1	61.5	44.9
------	------	------	------	------

CoME	35.2	64.5	62.5	18.6
------	------	------	------	------

CoME	<b>68.9</b>	<b>90.6</b>	<b>87.8</b>	<b>47.4</b>
------	-------------	-------------	-------------	-------------

We deeply appreciate your invaluable feedback and believe our response effectively addresses your concerns.

:ok: Finished ACLPUBcheck

:DF: [pdf](#)

I prefer a digital conference handbook

; have been developed to address such  
pose Conflict-free Model Editing (CoME), a  
e knowledge interference, allowing new  
ts, we demonstrate that CoME improves  
crucial for enhancing model editing

**editing. We are confident that the**

ns FT-based methods in Score, which  
prior knowledge. While this strategy can  
**demonstrates a more effective**

Handbook: 1

Final Check:

Paper Type:

Submission 1

Filter by review status

Everyone

## Paper Details

Decision: b

**Decision: A**

I prefer a digital conference handbook

I understand that no further revisions can be requested after the deadline.

Long

**Number:** 712

Reply type...   	Filter by author...   	Search keywords...	Sort: Newest First 	
None	Program Chairs	Submission712 Authors 		

## Decision

by Program Chairs  23 Jan 2025, 05:52 (modified: 23 Jan 2025, 07:50)  Program Chairs, Authors  Revisions

Accept to Main Conference

