

# Reasons and Solutions for the Decline in Model Performance after Editing

Xiusheng Huang<sup>1,2,3</sup>, Jiayang Liu<sup>1,2</sup>, Yequan Wang<sup>3</sup>, Kang Liu<sup>1,2</sup>

<sup>1</sup>The Key Laboratory of Cognition and Decision Intelligence for Complex Systems,  
Institute of Automation, Chinese Academy of Sciences

<sup>2</sup>School of Artificial Intelligence, University of Chinese Academy of Sciences

<sup>3</sup>Beijing Academy of Artificial Intelligence, Beijing, China



# Introduction

- Knowledge editing, while effective, may cause damage to the model:
  - > This study examines two factors causing this decline:
    - ~ Data perspective and Model perspective

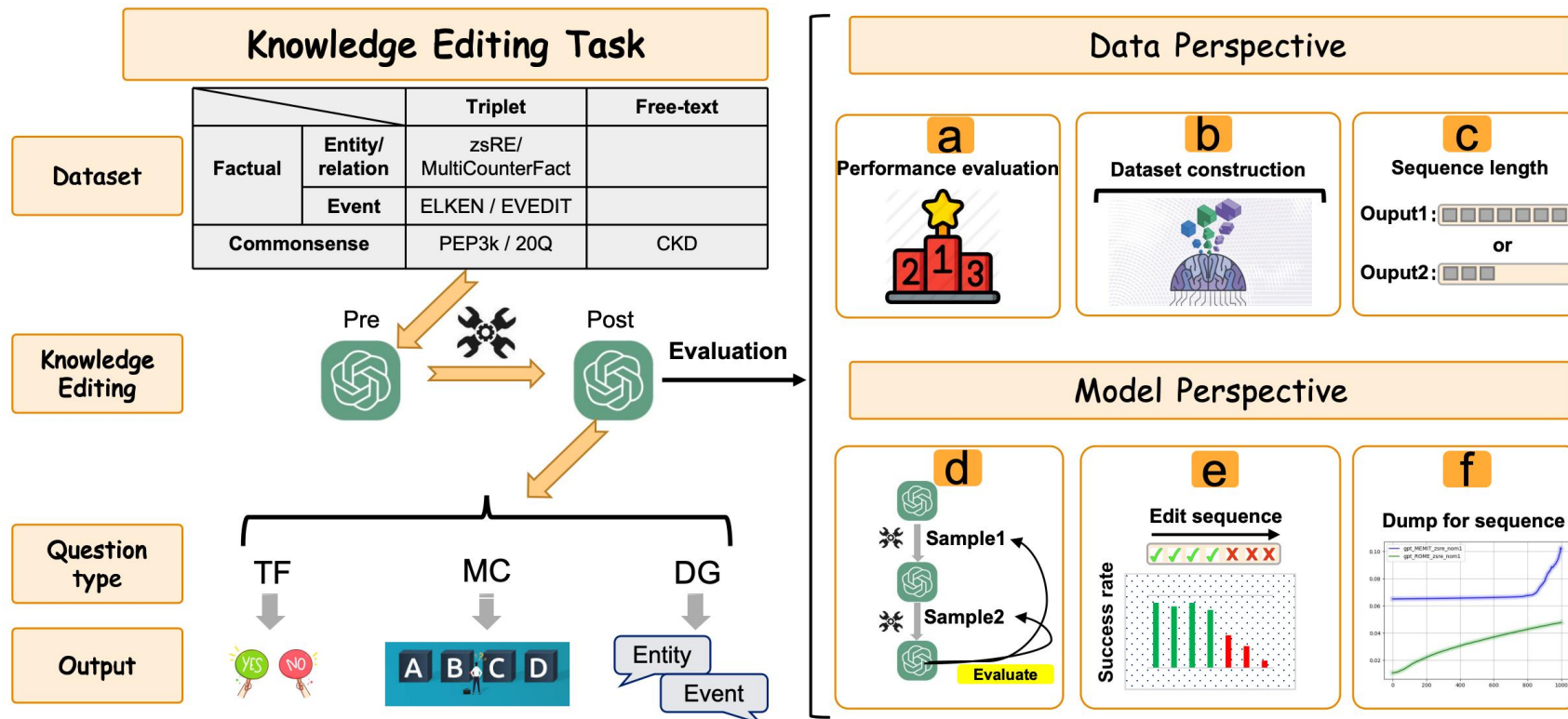


Fig 1. The framework of our work.

# Data perspective

- We conduct performance evaluation on multiple datasets:
  - Existing methods affect the performance of downstream tasks.

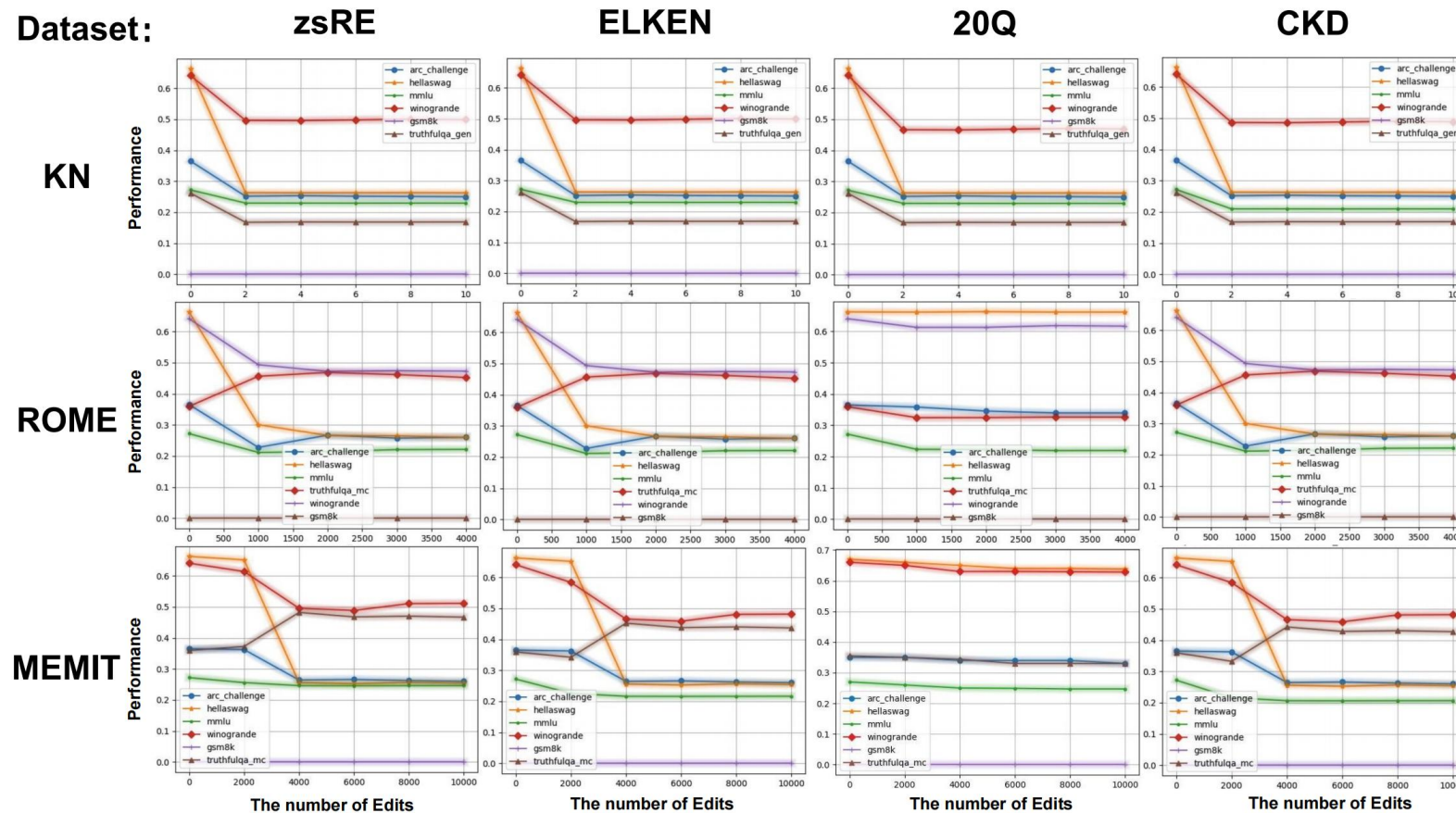


Fig 2. Evaluation results of different editing methods.

# Data perspective

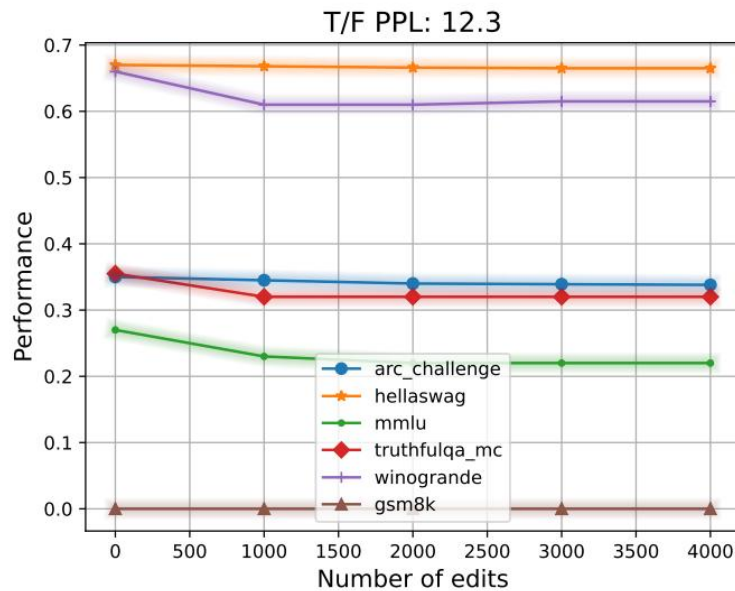
- To elucidate the impact of different editing objectives:
  - > We conduct: **MQD construction** -> impact assessment
    - ~ We created a **Multi-Question Dataset (MQD)** from ATOMIC
    - ~ Perplexity (PPL) values for the editing objectives: **297.4, 43.3, 12.3**
    - ~ Average length for the three question types: **23.44, 35.03, 13.38**

<b>ATOMIC Data Source: &lt; PersonX accepts PersonY appointment, as a result, PersonY shakes PersonX hand &gt;</b>		
<b>Category</b>	<b>Component Prompt</b>	<b>Target Answer</b>
DG	PersonX accepts PersonY appointment, resulting in PersonY	shakes PersonX hand
MQ	PersonX accepts PersonY appointment, resulting in PersonY ? Below are four options: (a) to give him a treat; (b) to forgive him; (c) to live happily ever after; (d) shakes PersonX hand. The correct option is	d
T/F	PersonX accepts PersonY appointment, resulting in PersonY shakes PersonX hand . Is this sentence logical? Please answer yes or no. A:	yes

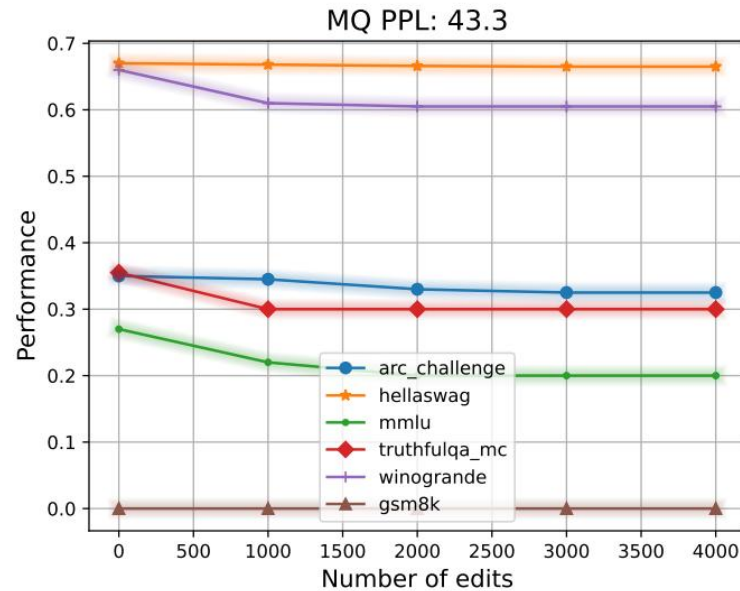
Table 1. An example of converting source data.

# Data perspective

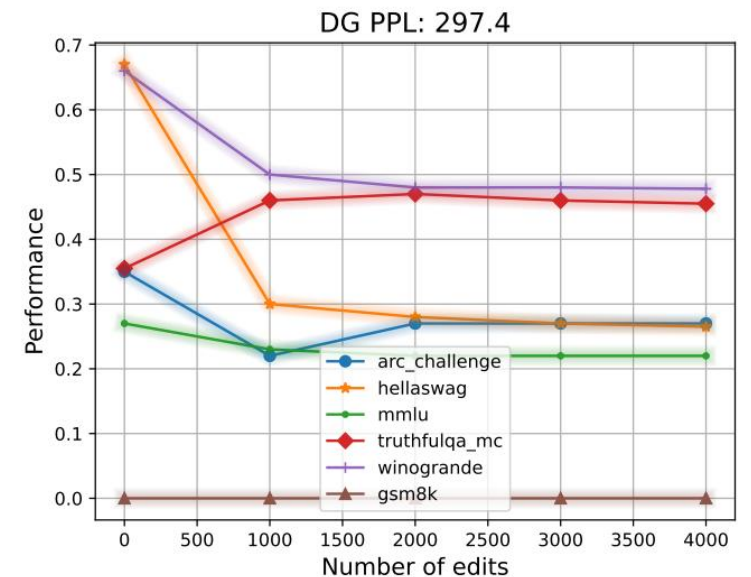
- To elucidate the impact of different editing objectives:
  - > We conduct: MQD construction -> **impact assessment**
    - ~ The decline in model performance after editing is attributed to the diversity of editing objectives and the length of tokens



(a) True/False questions (T/F).



(b) Multiple-choice (MQ).



(c) Directly generated (DG).

Fig 3. The performance of the model after editing.



# Model perspective

- To investigate the model factors contributing to degradation:
  - > We propose a new evaluation methods.

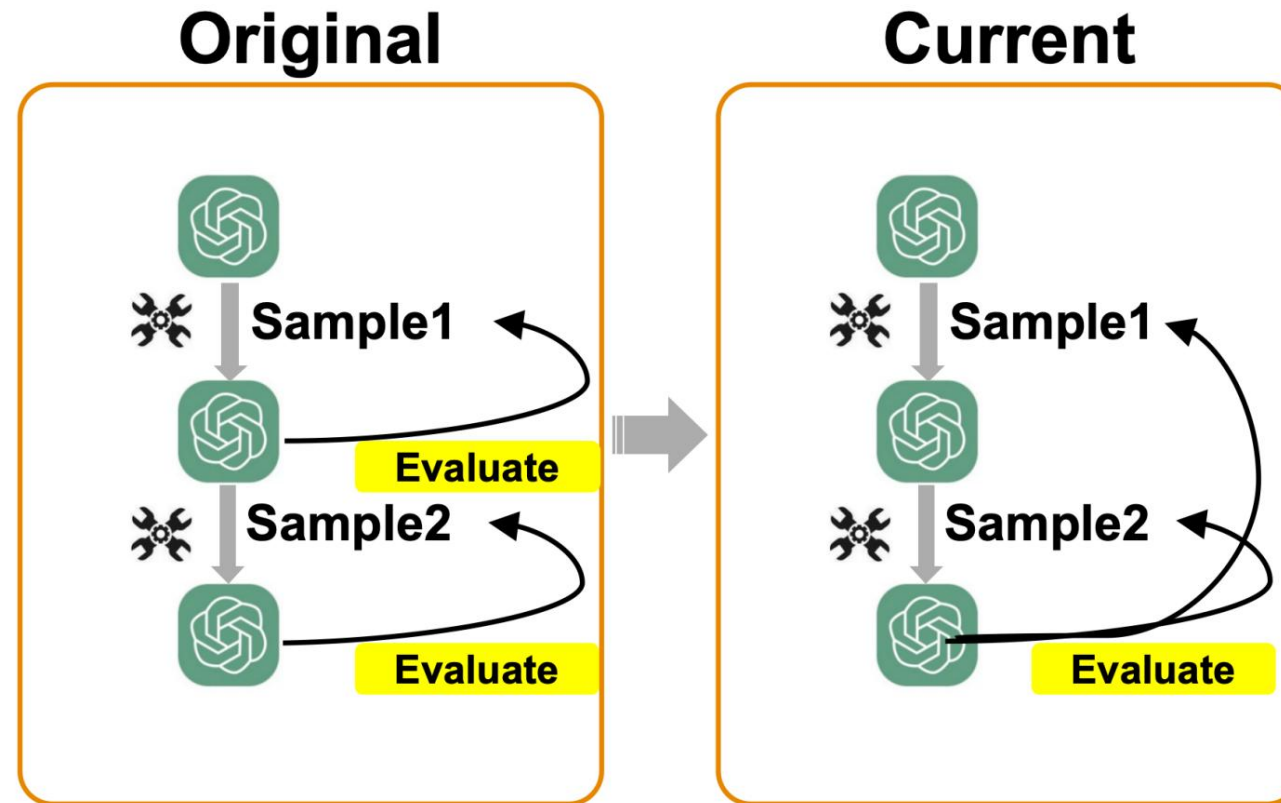


Fig 4. The original and current evaluation methods.

# Model perspective

- To investigate the model factors contributing to degradation:
  - > We evaluate the model's forgetfulness towards edited facts.

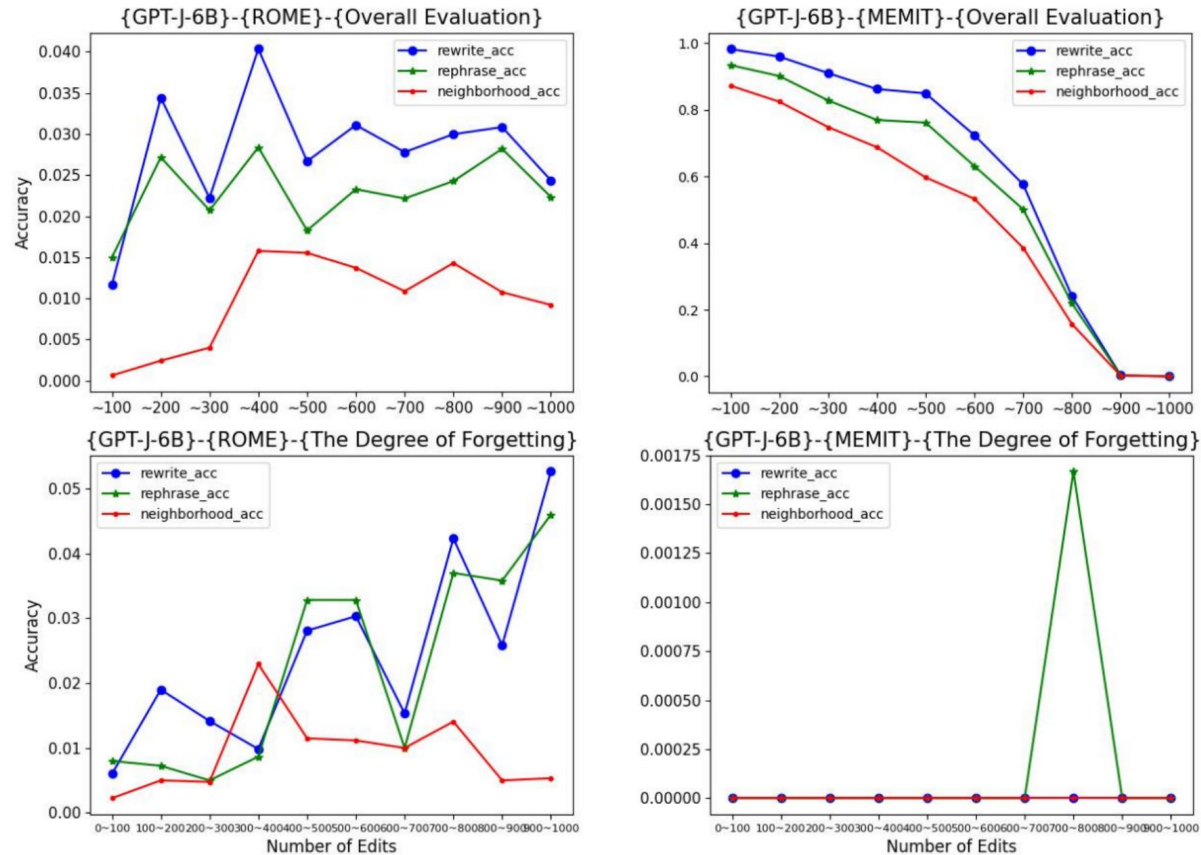
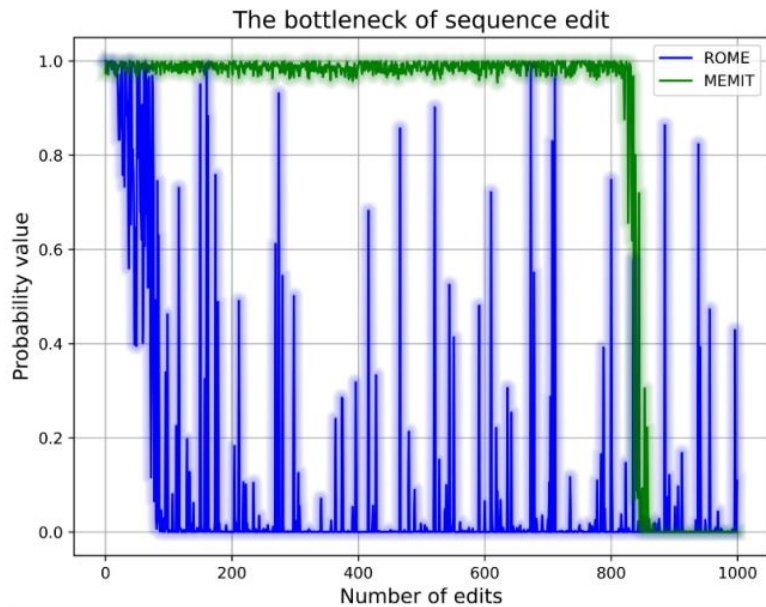


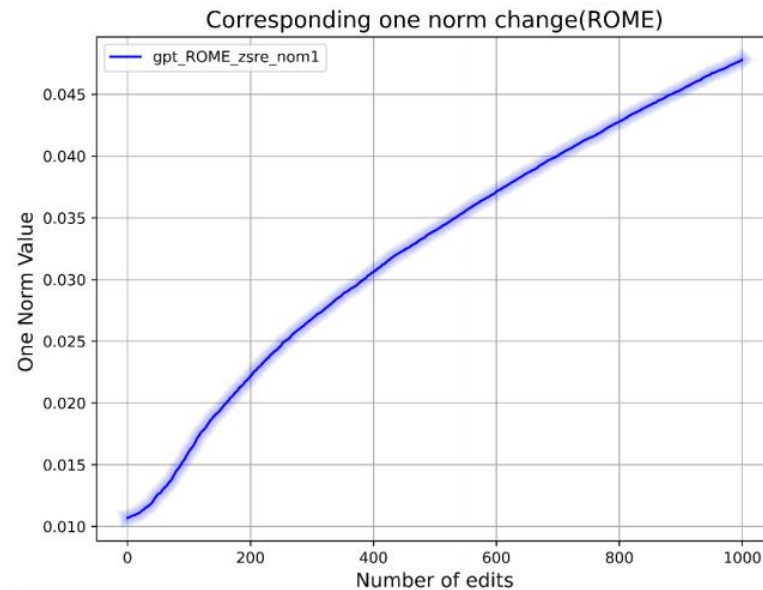
Fig 5. Assessment of forgetting ability of models.

# Model perspective

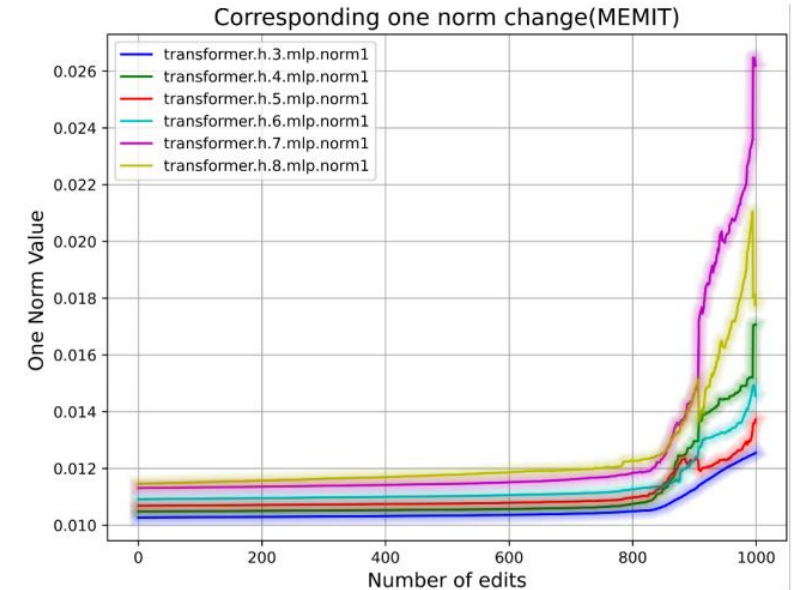
- To investigate the model factors contributing to degradation:
  - > Then, we examine the bottleneck of sequence edit.
    - ~ The decline in model performance after editing is due to the explosive growth of norms in the editing layers during the editing process



(a) The probability value.



(b) L1-norm on ROME.



(c) L1-norm on MEMIT.

Fig 6. The bottleneck and L1-norm correspondence in sequence editing.



# D4S method

- In order to mitigate the increase of the norm, we propose D4S:

> First, let's review the MEMIT process:

~ For the  $l^{th}$  FFN layer that needs to be modified:

$$k_i^l = \frac{1}{N} \sum_{k=1}^N \sigma(W_{in}^l \gamma(h_i^{l-1} [pre_k \oplus p(s_i)])), r_i^l = \frac{z_i - h_i^L}{L - l + 1}$$

$$\Delta^l = (R^l K^{lT})(Cov^l + K^l K^{lT})^{-1}, W_{out}^l \leftarrow W_{out}^l + \Delta^l$$

$$R^l = [r_1^l, \dots, r_n^l], K^l = [k_1^l, \dots, k_n^l], Cov^l = K_0^l K_0^{lT}$$

~ where  $k_i^l$  and  $r_i^l$  are the input and target output respectively. And  $K_0^l$  are the inputs of irrelevant knowledge

- In order to mitigate the increase of the norm, we propose D4S:

> For our Dump for Sequence (D4S) method:

~ We can split the  $\Delta^l$  into two linear segments:

$$R^l K^{lT} = [r_1^l, \dots, r_n^l][k_1^l, \dots, k_n^l]^T = \sum_{i=1}^n r_i^l k_i^{lT}$$

$$K^l K^{lT} = [k_1^l, \dots, k_n^l][k_1^l, \dots, k_n^l]^T = \sum_{i=1}^n k_i^l k_i^{lT}$$

~ For new knowledge to be edited:

$$R^l K^{lT'} = R^l K^{lT} + r_{n+1}^l k_{n+1}^{lT}$$

$$K^l K^{lT'} = K^l K^{lT} + k_{n+1}^l k_{n+1}^{lT}$$

# Experiments

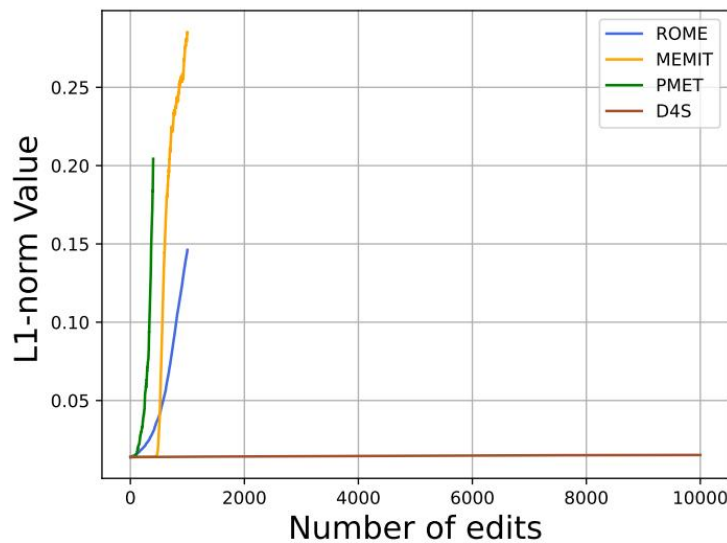
- The experimental results demonstrates the effectiveness of D4S:
  - > The experiments are designed to answer two research questions:
    - ~ How does D4S perform in sequence editing compared to other methods?
    - ~ How much damage does D4S do to the model?

Model	Method	<i>Edits</i> = 100				<i>Edits</i> = 500				<i>Edits</i> = 1000			
		<i>Eff.</i>	<i>Par.</i>	<i>Spe.</i>	<i>Avg.</i>	<i>Eff.</i>	<i>Par.</i>	<i>Spe.</i>	<i>Avg.</i>	<i>Eff.</i>	<i>Par.</i>	<i>Spe.</i>	<i>Avg.</i>
Llama	FT	30.47	25.78	12.19	22.81	19.43	18.06	6.56	14.68	16.13	13.89	5.96	12.00
	ROME	<u>97.68</u>	<b>92.39</b>	90.22	<b>93.43</b>	48.68	<u>47.72</u>	32.56	42.99	21.84	<u>20.74</u>	4.01	15.53
	MEMIT	92.99	83.22	<u>94.02</u>	90.08	2.90	2.90	2.51	2.77	2.85	2.85	2.74	2.82
	PMET	0.24	0.35	0.46	0.35	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
	GRACE	<b>99.33</b>	0.53	<b>99.88</b>	66.58	<b>98.87</b>	0.59	<b>99.72</b>	<u>66.39</u>	<b>98.94</b>	0.35	<b>99.77</b>	<u>66.35</u>
	<b>D4S(ours)</b>	95.53	<u>86.45</u>	93.44	<u>91.81</u>	<u>91.19</u>	<b>81.48</b>	<u>83.34</u>	<b>85.34</b>	<u>87.36</u>	<b>79.30</b>	<u>78.99</u>	<b>81.88</b>
GPT	FT	15.28	7.40	0.42	7.70	12.95	7.84	0.23	7.01	7.03	<u>4.15</u>	0.11	3.77
	ROME	1.17	1.50	0.06	0.91	2.67	1.83	1.55	2.02	2.44	2.23	0.92	1.86
	MEMIT	<u>98.25</u>	<u>93.41</u>	<u>87.24</u>	<u>92.97</u>	84.94	<u>76.15</u>	59.72	73.60	0.00	0.02	0.00	0.01
	PMET	0.33	0.33	0.00	0.22	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
	GRACE	<b>100.00</b>	0.40	<b>100.00</b>	66.8	<b>100.00</b>	0.15	<b>100.00</b>	66.72	<b>100.00</b>	0.07	<b>100.00</b>	66.69
	<b>D4S(ours)</b>	97.72	<b>97.01</b>	85.30	<b>93.34</b>	<u>97.42</u>	<b>93.12</b>	<u>74.65</u>	<b>88.40</b>	<u>95.98</u>	<b>91.17</b>	<u>70.17</u>	<b>88.77</b>

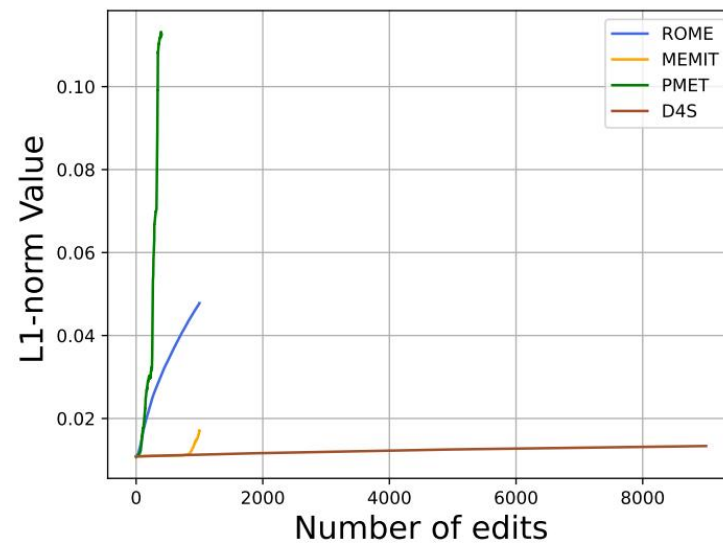
Table 2. Sequential edits on GPT and Llama with ZsRE dataset.

# Experiments

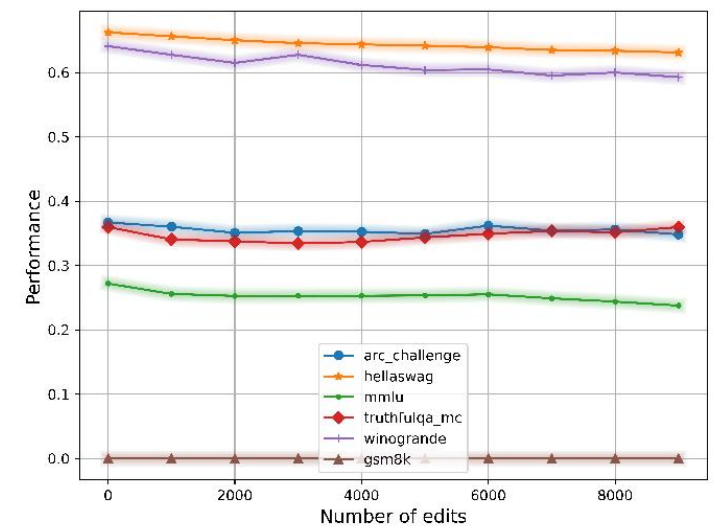
- The experimental results demonstrates the effectiveness of D4S:
  - > The experiments are designed to answer two research questions:
    - ~ How does D4S perform in sequence editing compared to other methods?
    - ~ **How much damage does D4S do to the model?**



(a) Llama\_L1-Norm



(b) GPT\_L1-Norm



(c) Downstream tasks

Fig 7. Norms of weight and performance of the edited model.

---

Thanks !