

## Knowledge Editing

1. Learning to Edit: Aligning LLMs with Knowledge Editing
2. Lifelong Knowledge Editing for LLMs with Retrieval-Augmented Continuous Prompt Learning

주간 세미나  
2025.01.09

# Knowledge Editing (지식 편집)이란

---

## [정의]

**Updating specific knowledge** stored in a **pre-trained language model**

- (1) **Preserving the overall structure and learned knowledge** of the existing model
- (2) **Selectively changing specific information** without broadly altering the model's behavior
- (3) **Efficiently and locally updating** the model's knowledge to address new facts or correct errors

## [목적]

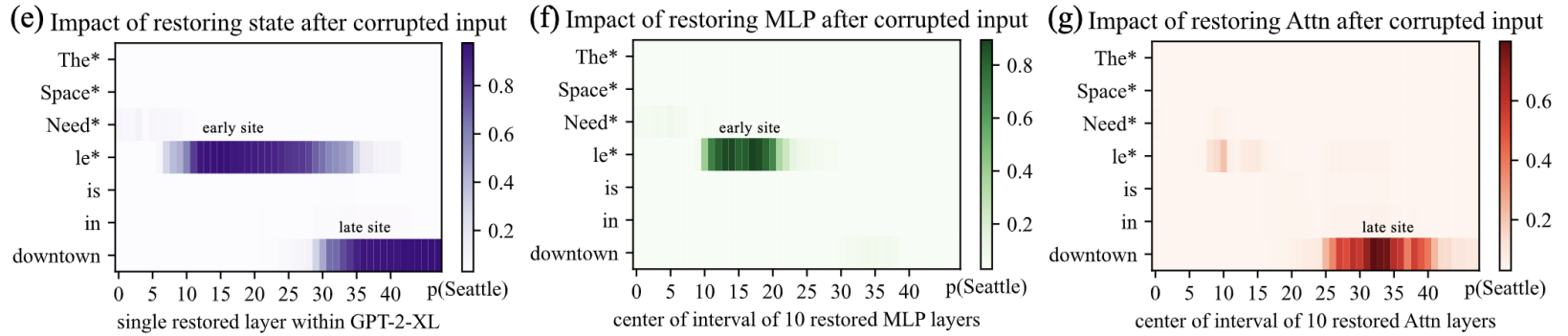
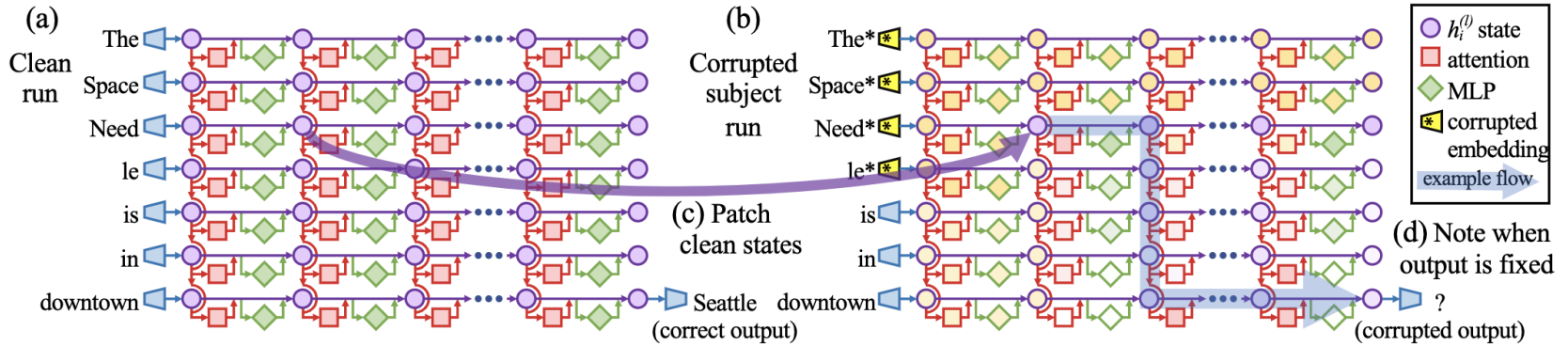
While retraining or finetuning can edit a model's predictions, doing this frequently is often too computationally expensive. LLaMA, for instance, was trained for 21 days on 2,048 A100 GPUs, costing over \$2.4M and emitting over 1,000 tons of CO<sub>2</sub>. – **GRACE 2023**

To address this need, model editing is an emerging area of research that aims to enable fast, data-efficient updates to a pre-trained base model's behavior for only a small region of the domain, without damaging model performance on other inputs of interest – **SERAC 2022**

# Knowledge Editing (지식 편집)이란

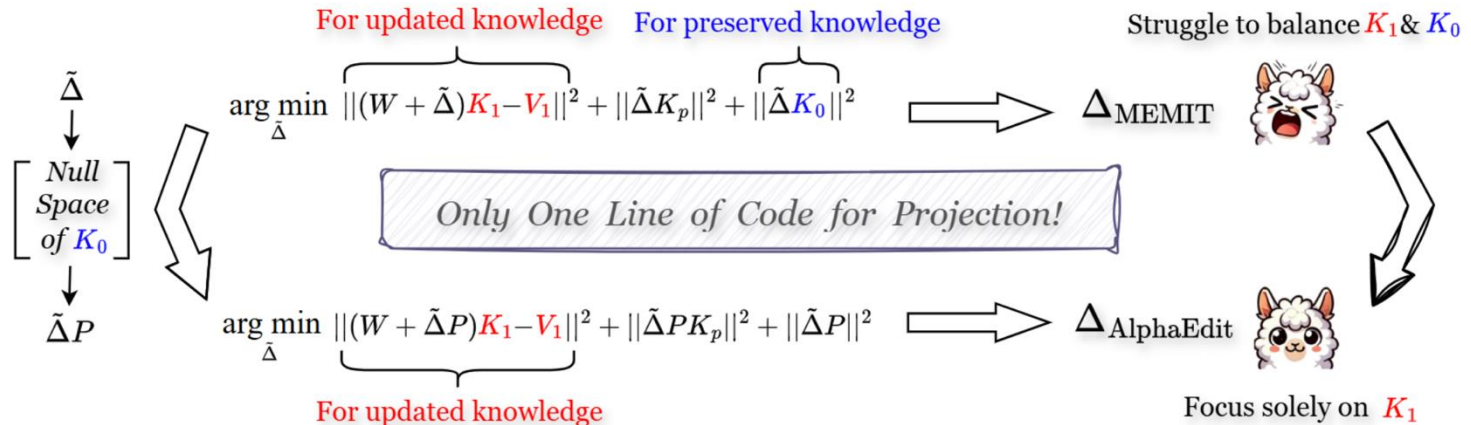
## Locate-then-Edit

### [ROME]



### [MEMIT]

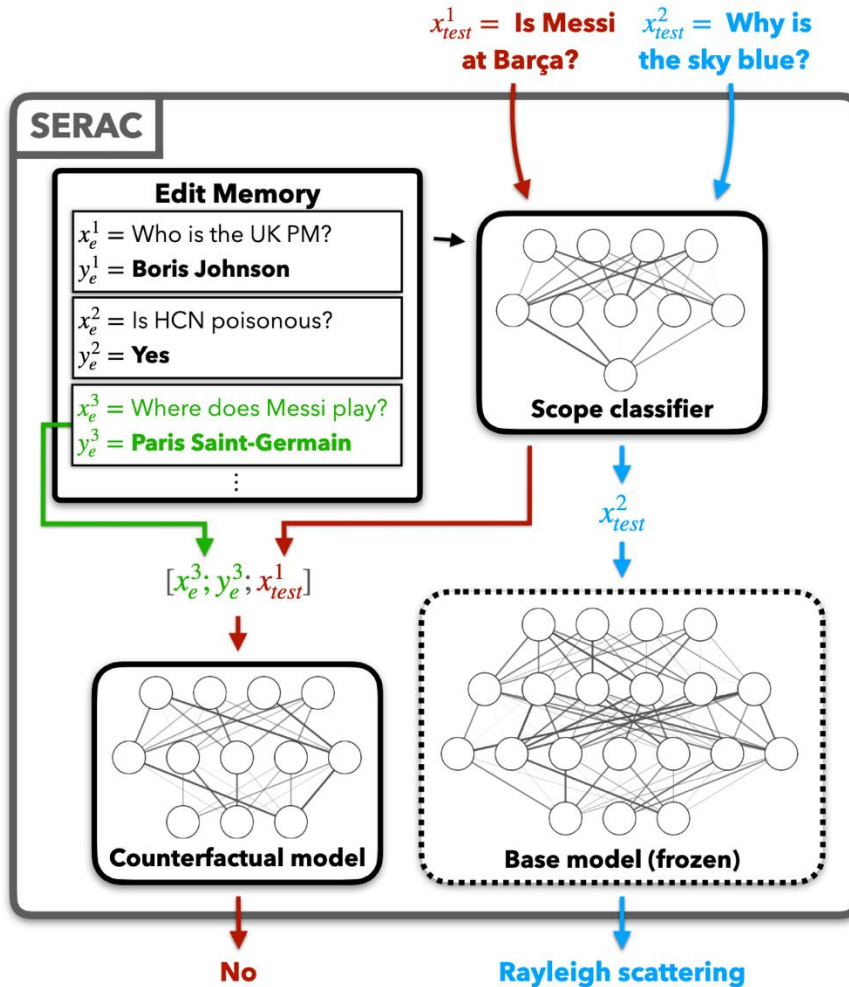
### [AlphaEdit]



# Knowledge Editing (지식 편집)이란

## Memory-based learning

[SERAC]



[GRACE]

**Algorithm 1:** Update Codebook at layer  $l$ .

**Input:**  $\mathcal{C} = \{(\mathbb{K}_i, \mathbb{V}_i, \epsilon_i)\}_{i=0}^{C-1}$ , codebook

**Input:**  $f(\cdot)$ , model

**Input:**  $y_t$ , desired label

**Input:**  $x_t$ , edit input for which  $f(x_t) \neq y_t$

**Input:**  $\epsilon_{\text{init}}$ , initial  $\epsilon$

**Input:**  $d(\cdot)$ , distance function

**Output:**  $\mathcal{C}$ , updated codebook

$C = \|\mathcal{C}\|$

$\hat{y}, h^{l-1} = f^L(x_t), f^{l-1}(x_t)$

$d_{\min}, i = \min_i (d(h^{l-1}, \mathbb{K}_i))$

If  $d_{\min} > \epsilon_i + \epsilon_{\text{init}}$  or  $C = 0$ :

#  $h^{l-1}$  far from existing entries or empty  $\mathcal{C}$

$v_{\text{new}} = \text{finetune on } P_f(y|v_{\text{init}})$

$\mathcal{C}_C = (h^{l-1}, v_{\text{new}}, \epsilon_{\text{init}})$  # Add entry

Else:

#  $h^{l-1}$  near existing entries

If  $f^L(k_i) = y$ :

# Same label  $\rightarrow$  Expand

$\mathcal{C}_i := (k_i, v_i, \epsilon_i + \epsilon_{\text{init}})$

Else:

# Different label  $\rightarrow$  Split

$\mathcal{C}_i = (k_i, v_i, d_{\min}/2)$  # Update entry  $i$

$v_{\text{new}} = \text{finetune on } P_f(y|v_{\text{init}})$

$\mathcal{C}_C = (h^{l-1}, v_{\text{new}}, d_{\min}/2)$  # Add entry

**return:**  $\mathcal{C}$

# Knowledge Editing (지식 편집)이란

## Task Formation

The objective of knowledge editing is to efficiently adjust the behavior of an initial base LLM  $f_\theta$ , where  $\theta$  represents the model's parameters, in response to specific **edit descriptors**  $\{(x_i^*, y_i^*)\}_{i \in [1, N]}$

$x_i^*$  : edit input trigger, query, question, key...

$y_i^*$  : edit target, target answer, desired label, value...

$N$  signifies the total number of edit descriptors.

The efficacy of knowledge editing is evaluated among several dimensions:

(1) **Reliability (Edit Success)**: measures the average accuracy of the post-edit model  $f_\theta^*$

$$\mathbb{E}_{(x_i^*, y_i^*)} \mathbb{1} \left\{ \arg \max_y f_\theta^*(y | x_i^*) = y_i^* \right\}$$

(2) **Portability**: evaluates how well updated knowledge transfers to related queries, enhancing the model's utility in varied contexts

(3) **Locality**: assesses the precision of edits, ensuring modifications are confined to targeted areas without affecting unrelated knowledge

(4) **Fluency**: quantifies the linguistic quality of the model's output post-edit, focusing on coherence and diversity to avoid repetitive pattern

# 1. Learning to Edit: Aligning LLMs with Knowledge Editing

**Yuxin Jiang<sup>1,2\*</sup>, Yufei Wang<sup>3</sup>, Chuhan Wu<sup>3</sup>, Wanjun Zhong<sup>3</sup>,  
Xingshan Zeng<sup>3</sup>, Jiahui Gao<sup>3</sup>, Liangyou Li<sup>3</sup>, Xin Jiang<sup>3</sup>,  
Lifeng Shang<sup>3</sup>, Ruiming Tang<sup>3</sup>, Qun Liu<sup>3</sup>, Wei Wang<sup>1,2</sup>**

The Hong Kong University of Science and Technology (Guangzhou)<sup>1</sup>,  
The Hong Kong University of Science and Technology<sup>2</sup>, Huawei Noah's Ark Lab<sup>3</sup>  
yjiangcm@connect.ust.hk, wang.yufei1@huawei.com, weiwcs@ust.hk

**ACL 2024**

# Introduction

---

## Previous Knowledge Editing Approaches..

(1) Rely on **auxiliary modules** or models to either predict the **LLM's weight adjustments**

→ Meta-learning 기반의 방식들의 한계

(2) Function as **scope classifiers** for query response applicability

→ Memory 기반의 SERAC의 한계

(3) Localization results from **Causal Tracing are statistically uncorrelated** with the success of an edit injecting a new fact into MLP weights.

→ Locate-then-edit 기반의 한계

***“Teach a man to fish, and you feed him for a lifetime”***

→ To elicit LLMs' capabilities of following knowledge editing instructions, thereby empowering them to **effectively leverage the updated knowledge to answer the queries**

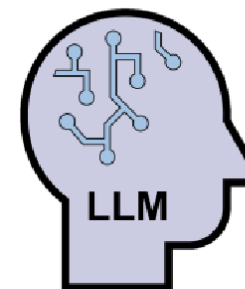
# Introduction

## Learning To Edit (LTE)

- To elicit LLMs' capabilities of following knowledge editing instructions, thereby empowering them to **effectively leverage the updated knowledge to answer the queries.**

### Previous Knowledge Editing Methods

The current British Prime Minister is Rishi Sunak



Who is married to the PM of the UK?



### Our Proposed LTE Framework

The current British Prime Minister is Rishi Sunak



Who is married to the PM of the UK?

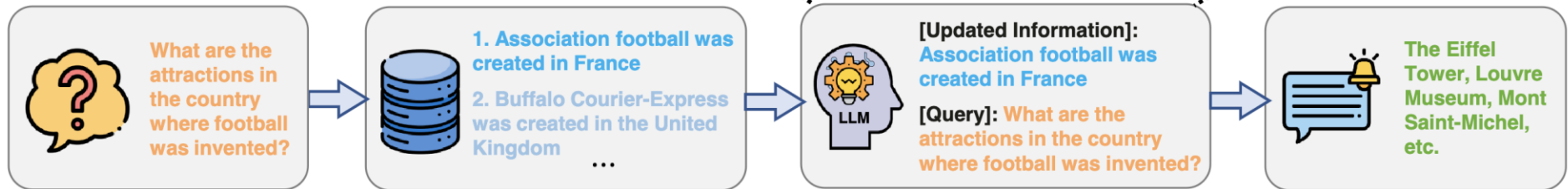
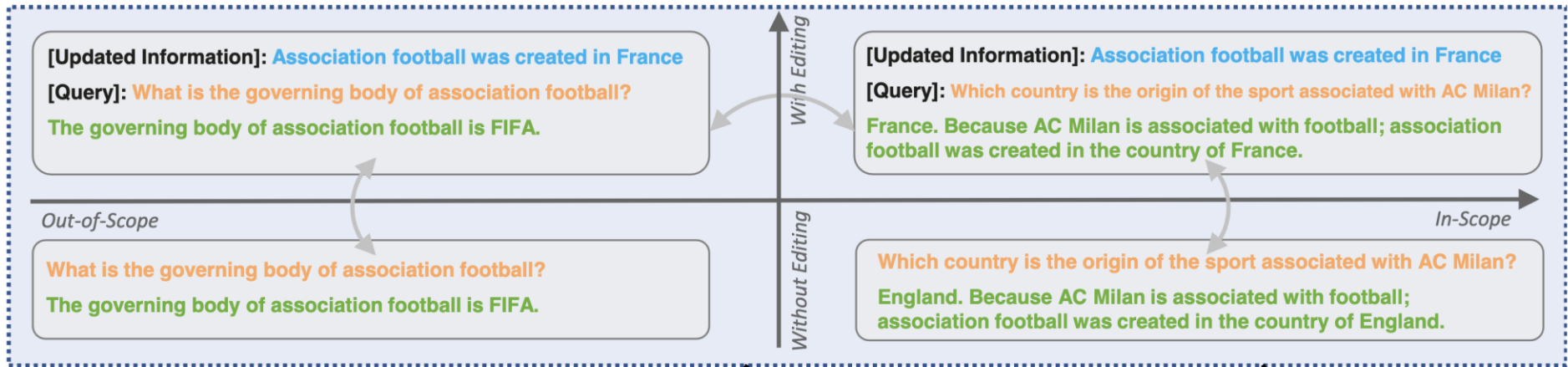


# Method

## Learning To Edit (LTE)

### Alignment Phase: Learning to Edit

■ Edit Descriptor ■ Query ■ Answer



### Inference Phase: On-the-fly Edit

**Learning to Edit (LTE)** framework to align LLMs with ever-changing, complicated, and diverse knowledge editing requests in real-time

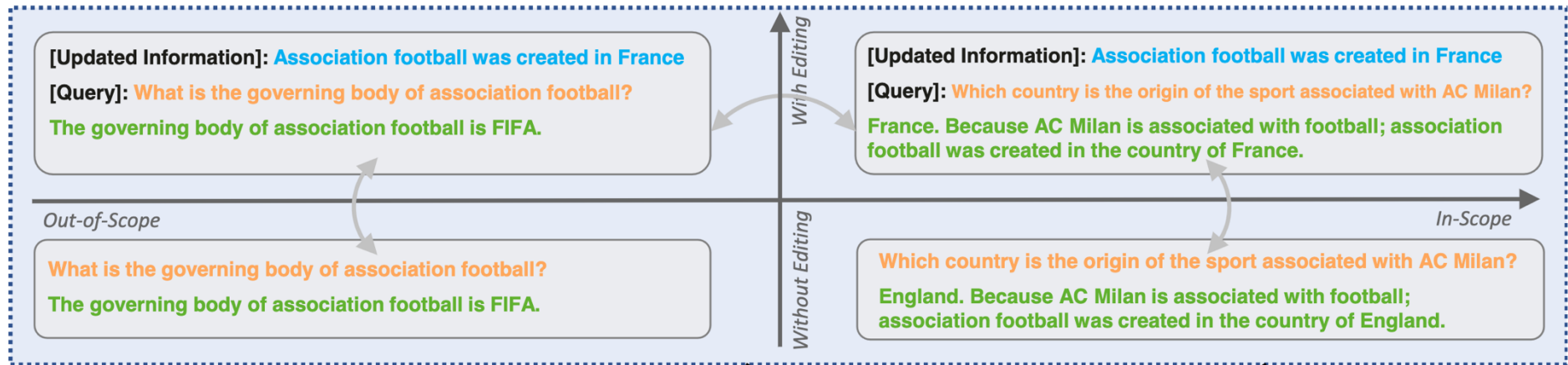
- (i) **Alignment Phase:** Knowledge Editing Prompt로 지식 편집이 가능하도록 학습하는 과정
- (ii) **Inference Phase:** 저장된 메모리 → 업데이트할 관련 지식 + Query로 Knowledge Editing

# Method

## Learning To Edit (LTE)

### Alignment Phase: Learning to Edit

■ Edit Descriptor ■ Query ■ Answer



### (i) Alignment Phase: Learning to Edit

“*[Updated Information]{edit descriptor}\n[Query]{query}*”

An **optimal knowledge editing method** must seamlessly **integrate new knowledge into the relevant content** within its edit scope, while **ensuring** the accuracy and integrity of information **outside this domain**

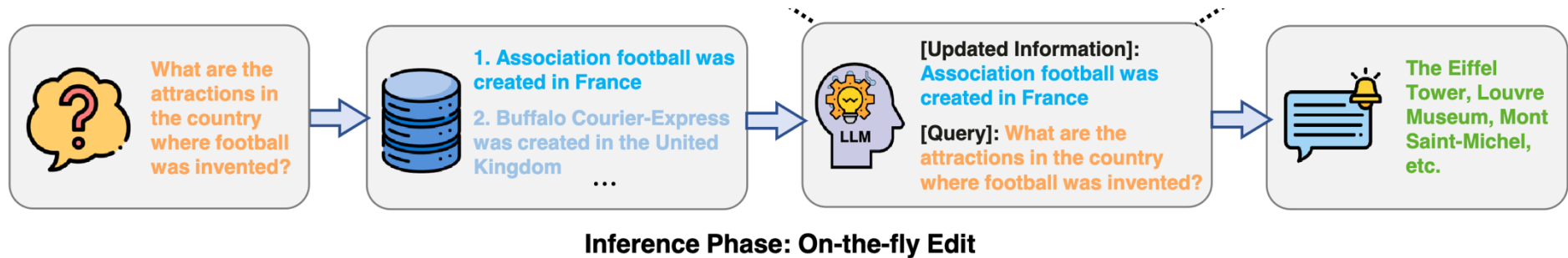
**(1) In-Scope Capability:** It also covers subject aliasing, ensuring the editing of one subject should not vary from its expression

**(2) Out-Scope Capability:** To maintain the integrity of unrelated attributes of the subject, ensuring no unintended alteration

**(3) Linguistic Capability:** Editing should not hinder the model’s proficiency in unrelated areas

# Method

## Learning To Edit (LTE)



### (ii) Inference Phase On-the-fly Edit

RAG-based mechanisms support real-time knowledge editing, retrieving the top  $k$  ( $k=3$ ) relevant updated knowledge from within stored memory (embed all the edit descriptors and create a vector memory) and incorporating it into model responses.

#### *threefold strategy*

규칙 1. Firstly, in 50% of cases, we directly use the exact edit descriptor

규칙 2. Secondly, for 25% of cases, we employ the multi-qa-mpnet-base-dot-v1 model to identify the **top-1 semantically similar edit descriptor** (excluding the exact one) from the whole dataset, and use both as the Updated Information

규칙 3. Lastly, for the remaining 25%, we retrieve the **top 2 semantically similar descriptors**, excluding the exact one, using all three as the Updated Information

# Experiments

## Main Results

Model	Dataset	Metric	SERAC	ICE	MEND	ROME	MEMIT	FT-L	FT	LTE	LTE-LoRA
LLaMA2-Chat-7B	ZsRE	Edit Succ.	<u>99.67</u>	66.01	96.74	96.57	83.07	54.65	36.88	<b>99.91</b>	<b>99.91</b>
		Portability	56.48	63.94	60.41	52.20	51.43	45.02	8.72	<u>78.98</u>	<b>79.63</b>
		Locality	30.23	23.14	<b>92.79</b>	27.14	25.46	71.12	0.31	<u>71.78</u>	67.99
		Fluency	410.89	541.14	524.33	<u>570.47</u>	559.72	474.18	471.29	<b>583.70</b>	544.52
	WikiBio	Edit Succ.	99.69	95.53	93.66	95.05	94.29	66.27	95.64	<b>99.87</b>	<u>99.76</u>
		Locality	69.79	47.90	69.51	46.96	51.56	60.14	13.38	<b>80.27</b>	<u>72.31</u>
		Fluency	606.95	<b>632.92</b>	609.39	<u>617.25</u>	616.65	604.00	589.22	614.26	611.94
	Recent	Edit Succ.	98.68	60.74	76.88	85.08	85.32	71.18	31.24	<b>99.99</b>	<u>99.97</u>
		Portability	63.52	36.93	50.11	37.45	37.94	48.71	15.91	<b>91.51</b>	<u>81.87</u>
		Locality	<b>100.00</b>	33.34	<u>92.87</u>	66.20	64.78	63.70	3.65	85.67	82.72
		Fluency	553.19	531.01	<u>586.34</u>	574.28	566.66	549.35	428.67	<b>586.76</b>	570.64
	Counterfact	Edit Succ.	<u>99.99</u>	69.83	78.82	83.21	83.41	51.12	26.78	<b>100.00</b>	99.97
		Portability	76.07	45.32	57.53	38.69	40.09	39.07	16.94	<b>89.69</b>	85.74
		Locality	<b>98.96</b>	32.38	<u>94.16</u>	65.40	63.68	62.51	0.29	84.76	85.11
		Fluency	549.91	547.22	<u>588.94</u>	578.84	568.58	544.80	483.71	<b>589.69</b>	574.14
	Average	Edit Succ.	99.51	73.03	86.53	89.98	86.52	60.81	47.64	<b>99.94</b>	<u>99.90</u>
		Portability	65.36	48.73	56.02	42.78	43.15	44.27	13.86	<b>86.73</b>	<u>82.41</u>
		Locality	74.75	34.19	<b>87.33</b>	51.43	51.37	64.37	4.41	<u>80.62</u>	77.03
		Fluency	530.24	563.07	577.25	<u>585.21</u>	577.90	543.08	493.22	<b>593.60</b>	575.31
	Qwen-Chat-7B	ZsRE	Edit Succ.	98.43	70.29	99.40	<b>99.90</b>	97.25	37.81	25.33	<u>99.72</u>
Portability			56.69	67.52	59.98	46.76	44.31	41.85	7.70	<b>82.92</b>	<u>80.16</u>
Locality			41.28	73.45	80.83	48.90	60.26	<b>87.70</b>	3.29	<u>80.99</u>	78.28
Fluency			495.12	556.86	544.07	562.88	<u>578.73</u>	557.86	538.10	<b>580.01</b>	543.35
WikiBio		Edit Succ.	99.39	94.60	93.38	98.79	96.10	60.19	34.63	<b>99.80</b>	<u>99.75</u>
		Locality	71.50	58.15	65.47	41.78	65.65	<b>80.41</b>	22.45	79.63	<u>80.34</u>
		Fluency	598.11	614.22	610.92	604.81	<u>623.49</u>	595.56	572.59	<b>634.73</b>	620.05
Recent		Edit Succ.	99.58	83.86	82.39	99.67	98.96	60.07	29.74	<b>99.73</b>	<u>99.68</u>
		Portability	67.22	58.24	57.92	50.84	49.38	42.02	14.33	<b>89.73</b>	<u>87.40</u>
		Locality	<b>100.00</b>	61.83	89.11	51.78	60.72	84.83	4.27	<u>89.25</u>	83.77
		Fluency	561.32	559.46	<u>610.72</u>	600.70	600.39	598.32	456.99	<b>615.59</b>	587.90
Counterfact		Edit Succ.	99.06	80.28	88.04	<b>99.44</b>	95.05	24.55	15.42	99.28	<u>99.35</u>
		Portability	79.28	53.80	52.99	40.63	34.50	20.14	11.38	<b>86.79</b>	85.33
		Locality	<u>92.70</u>	63.86	91.05	39.22	50.14	<b>92.74</b>	30.04	86.87	85.20
		Fluency	<u>568.05</u>	559.46	<u>619.87</u>	603.21	604.47	608.47	563.70	<b>622.91</b>	593.51
Average		Edit Succ.	99.12	82.26	90.80	99.45	96.84	45.66	26.28	<b>99.63</b>	<u>99.59</u>
		Portability	67.99	59.85	56.96	46.08	42.73	34.67	11.14	<b>86.48</b>	<u>84.30</u>
		Locality	76.37	64.32	81.62	45.42	59.19	<b>86.42</b>	15.01	<u>84.19</u>	81.90
		Fluency	555.65	572.50	596.40	592.90	<u>601.77</u>	590.05	532.85	<b>613.31</b>	586.20

# Experiments

## Mass & Sequential Editing

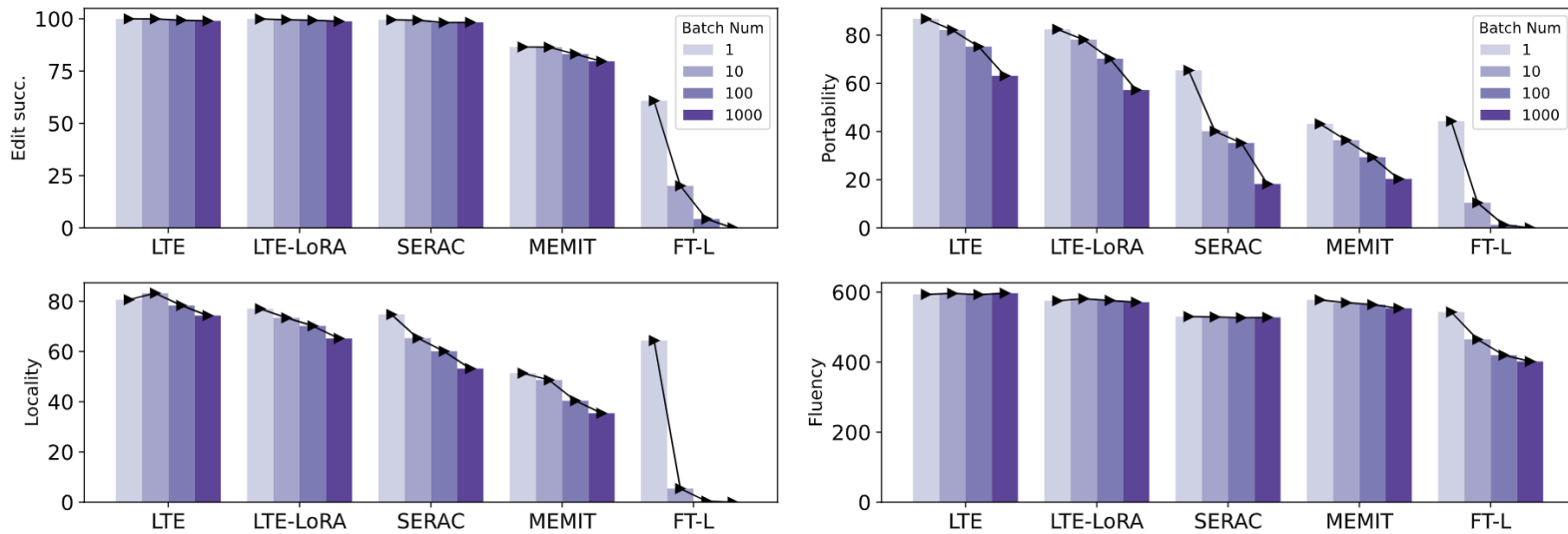
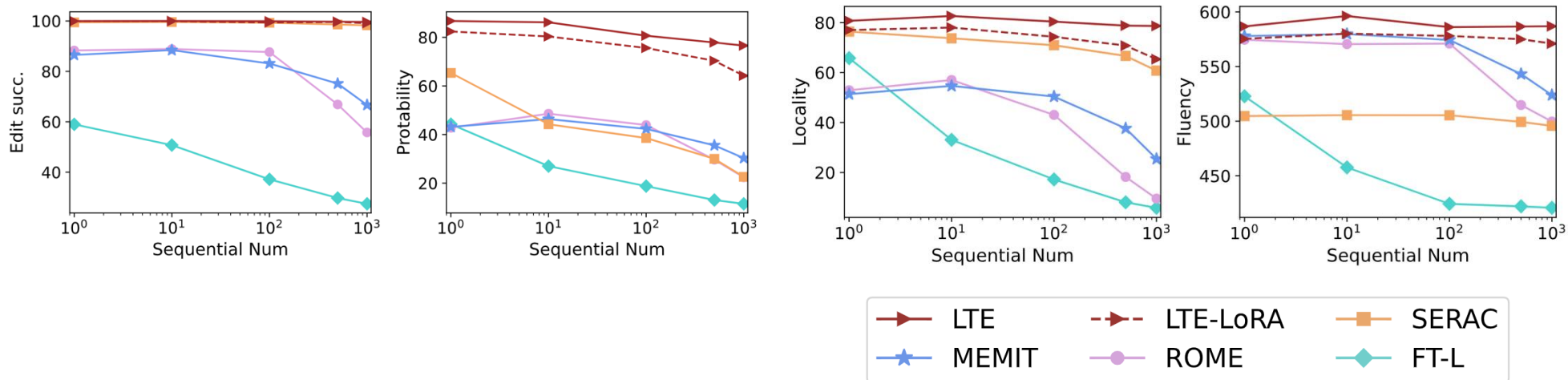


Figure 3: Averaged **Batch Editing** performance on four benchmarks against batch numbers in [1, 10, 100, 1000].



# Experiments

## Linguistic Capability

	CommonSenseQA	PIQA	XSum	MMLU	AGIEval	AlpacaEval	Average
<i>LLaMA2-Chat-7B</i>	<b>69.9</b>	<b>65.0</b>	22.3	40.4	26.1	71.4	49.2
LTE w/o editing	67.2	61.3	<b>22.4</b>	46.4	<b>26.5</b>	<b>73.3</b>	<b>49.5</b>
LTE w/ editing	67.1	62.6	<b>22.4</b>	<b>47.8</b>	23.8	71.6	49.2
<i>Qwen-Chat-7B</i>	<b>77.6</b>	<b>72.1</b>	28.8	56.6	41.3	77.8	59.0
LTE w/o editing	74.7	69.3	29.9	<b>59.3</b>	<b>41.9</b>	<b>79.2</b>	<b>59.1</b>
LTE w/ editing	75.3	70.0	<b>30.1</b>	58.2	40.7	78.4	58.8

Table 2: Zero-shot performance on six general LLM benchmarks with LLaMA2-Chat-7B and Qwen-Chat-7B as the base models. “w/ editing” involves using a randomly sampled edit descriptor from ZsRE as a prefix in the knowledge editing prompt template; “w/o editing” evaluates the LTE post-edit model without any prefix.

## Ablation

	S	P	L	F	G
LTE	99.94	86.73	80.62	593.60	49.5
-w/o in-scope training	<b>77.53</b>	<b>56.26</b>	80.72	589.04	49.0
-w/o out-of-scope training	99.92	86.89	<b>65.50</b>	592.66	49.2
-w/o free-text QA training	99.93	86.30	80.91	<b>587.75</b>	<b>43.9</b>
-w/o threefold strategy	99.78	86.51	80.22	593.40	49.5
-w/o training	<b>75.04</b>	<b>54.23</b>	<b>48.19</b>	592.73	49.2

	Seq_Num	Edit Succ.	Portability	Locality
LTE w/ 420M <i>R</i> top <i>k</i> = 3	10	100.00	86.16	82.64
	100	99.90	80.66	80.38
	1000	99.64	76.59	78.67
LTE w/ 80M <i>R</i> top <i>k</i> = 3	10	100.00	83.38	78.65
	100	99.81	79.92	80.40
	1000	99.61	75.67	79.43
LTE w/ 420M <i>R</i> top <i>k</i> = 2	10	100.00	85.69	81.59
	100	99.85	80.05	80.67
	1000	99.63	76.27	78.05
LTE w/ 420M <i>R</i> top <i>k</i> = 1	10	100.00	84.01	81.96
	100	99.83	79.48	80.11
	1000	99.56	75.93	78.89

## 2. Lifelong Knowledge Editing for LLMs with Retrieval-Augmented Continuous Prompt Learning

**Qizhou Chen<sup>1,2\*</sup>, Taolin Zhang<sup>2\*</sup>, Xiaofeng He<sup>1,3</sup>, Dongyang Li<sup>1,2</sup>,  
Chengyu Wang<sup>4†</sup>, Longtao Huang<sup>2</sup>, Hui Xue<sup>2</sup>**

<sup>1</sup> East China Normal University, Shanghai, China <sup>2</sup> Alibaba Group, Hangzhou, China

<sup>3</sup>NPPA Key Laboratory of Publishing Integration Development, ECNUP, Shanghai, China

<sup>4</sup> Alibaba Cloud Computing, Hangzhou, China

chen\_qizhou@outlook.com, zhangt10519@gmail.com, chengyu.wcy@alibaba-inc.com

**EMNLP 2024**

# Introduction

## Previous Knowledge Editing Approaches..

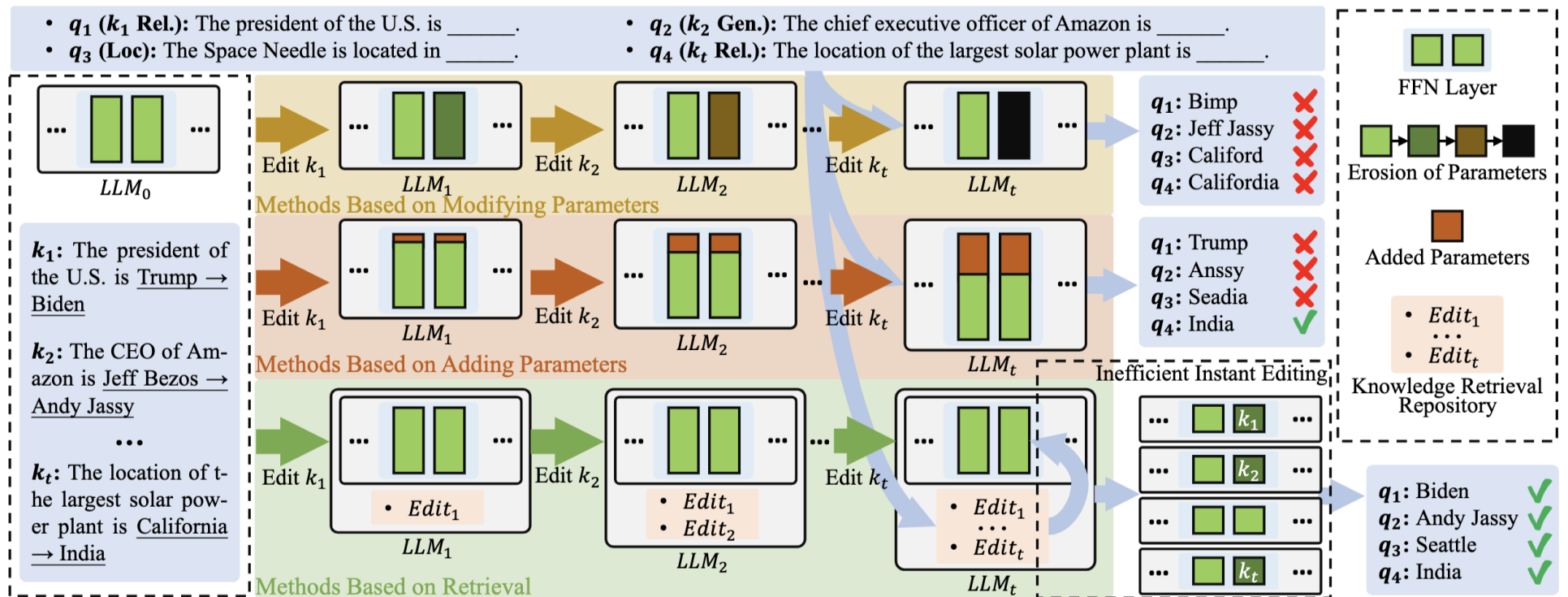


Figure 1: Comparison among three types of methods in lifelong editing scenarios. Modifying parameters and adding extra parameters result in the degradation of LLM performance as editing progresses. In contrast, retrieval-based editors store knowledge in a repository and apply knowledge editing on the fly, which maintains the LLM unchanged and relieves it from accumulating parameter offsets or adding extra parameters. (Best viewed in color)



# Method

## RECIPE (a RetriEval-augmented Continuous Prompt LEarning)

### 1. Knowledge Continuous Prompt Learning

1) Avoids the shortcomings of LTE

(1) Overly long editing prefixes can reduce model inference speed

(2) Full-parameter fine-tuning also increases the risk of overfitting

→ P-tuning (Li and Liang 2021; Liu et al., 2022) based continuous prompt learning

→ Using trainable word embedding vectors as prompts

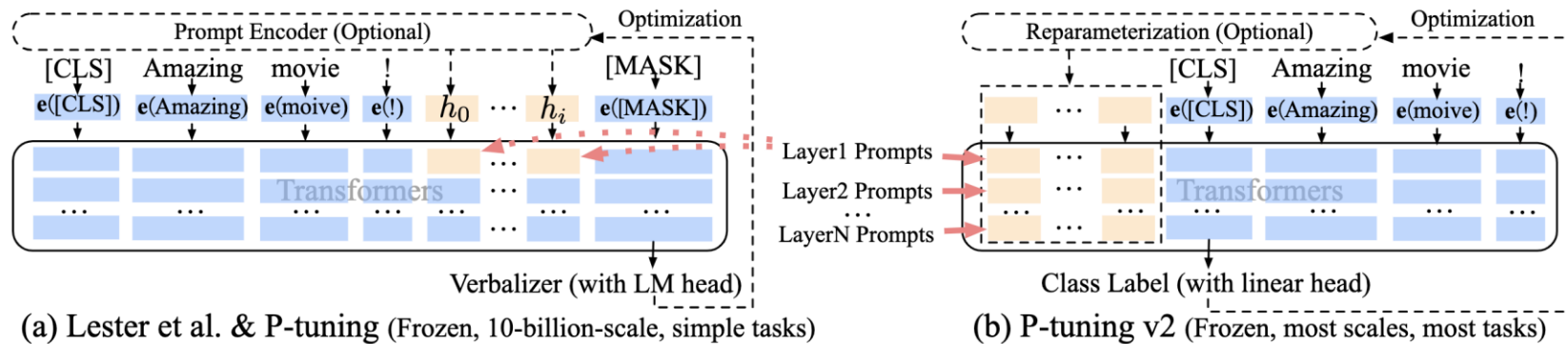


Figure 2: From Lester et al. (2021) & P-tuning to P-tuning v2. Orange blocks (i.e.,  $h_0, \dots, h_i$ ) refer to trainable prompt embeddings; blue blocks are embeddings stored or computed by frozen pre-trained language models.

## RECIPE (a RetriEval-augmented Continuous Prompt LEarning)

### 2. Dynamic Prompt Retrieval with Knowledge Sentinel

1) Mapping knowledge statements and queries into the same representational space

(1) 일반적으로는 setting a fixed similarity threshold를 설정함

→ **Dynamic thresholds**

→ **Knowledge Sentinel (KS)** (a trainable embedding representation)

→ Joint training with the prompt encoder and KS module

# Method

## RECIPE (a RetriEval-augmented Continuous Prompt lEarning)

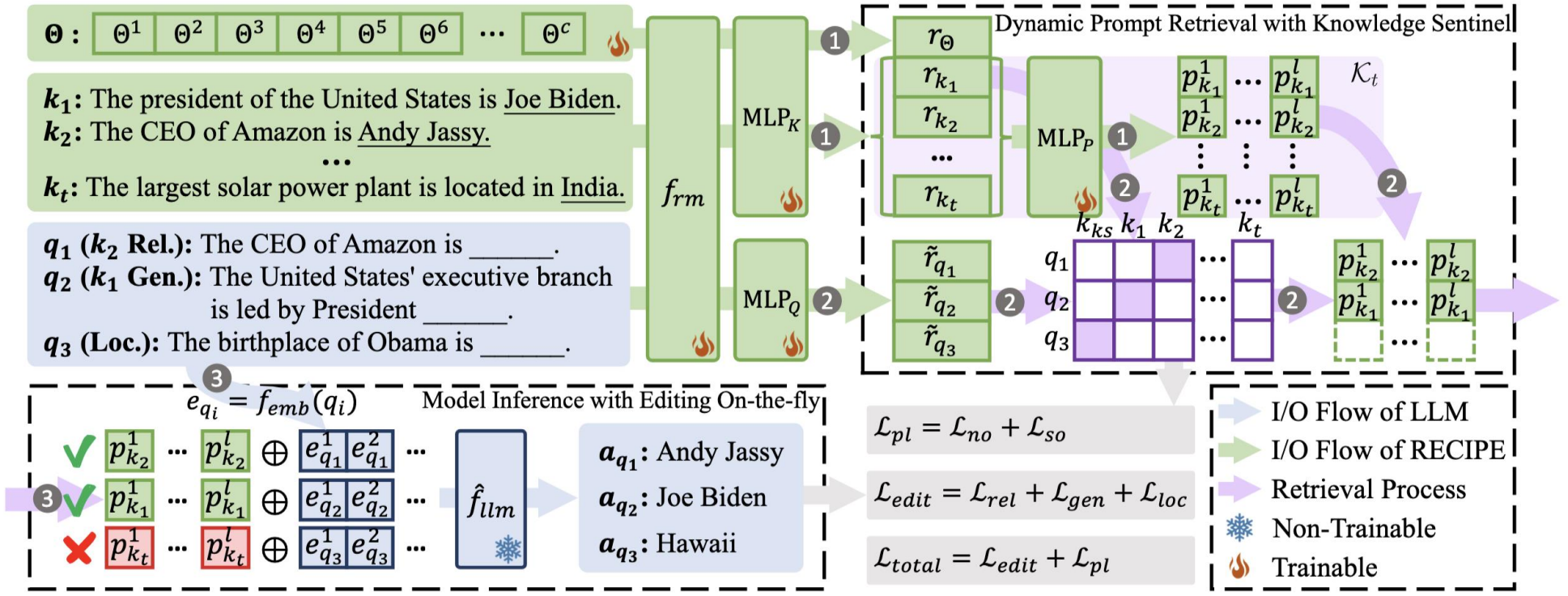
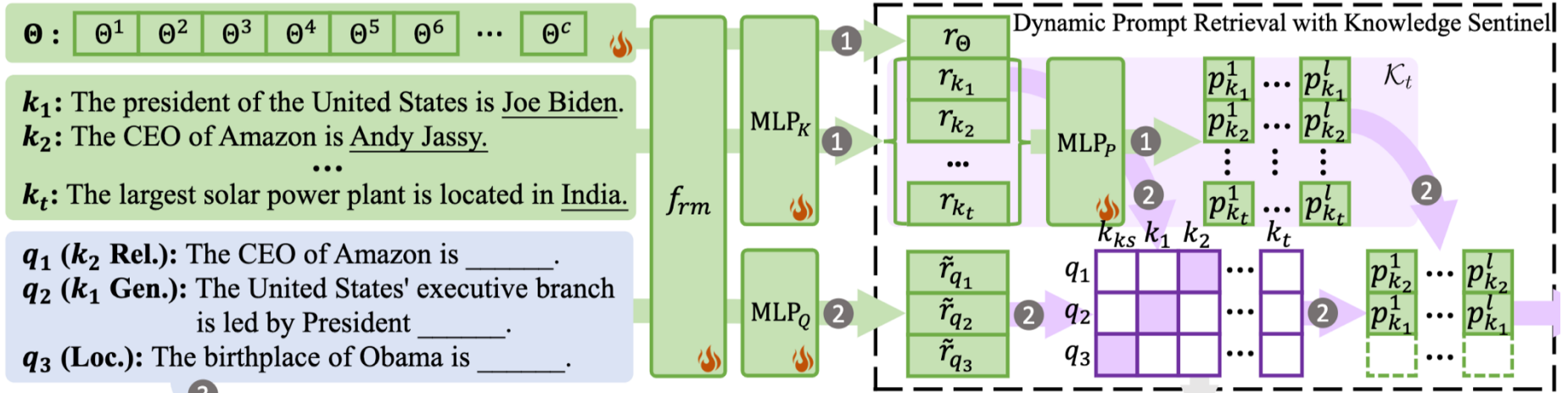


Figure 2: Illustration of the RECIPE framework. Process 1 constructs and updates the knowledge retrieval repository  $\mathcal{K}_t$ . During the inference stage, Process 2 retrieves query-related prompts from  $\mathcal{K}_t$ . Process 3 utilizes the retrieved continuous prompts to correct the LLM’s response. For lifelong editing, the repository can be continuously updated (e.g., from  $\mathcal{K}_{t-1}$  to  $\mathcal{K}_t$ ) with each new insertion of knowledge and prompts.

# Method

## RECIPE (a RetriEval-augmented Continuous Prompt Learning)

### 1. Construction and Update of Knowledge Retrieval Repository



1) Knowledge retrieval repository:  $K_0 = \{\}$

(1)  $K_{t-1}$  to  $K_t$  by adding a new key-value pair at each timestep  $t$

Given a new knowledge statement  $K_t$ ,

the knowledge representation is achieved through an encoder  $f_{rm}$  (RoBERTa)  $r_{k_t} = \mathbf{MLP}_K(f_{rm}(k_t))$

the continuous prompt  $p_{k_t}$  is generated through another MLP  $p_{k_t} = f_{resp}(\mathbf{MLP}_P(r_{k_t}))$

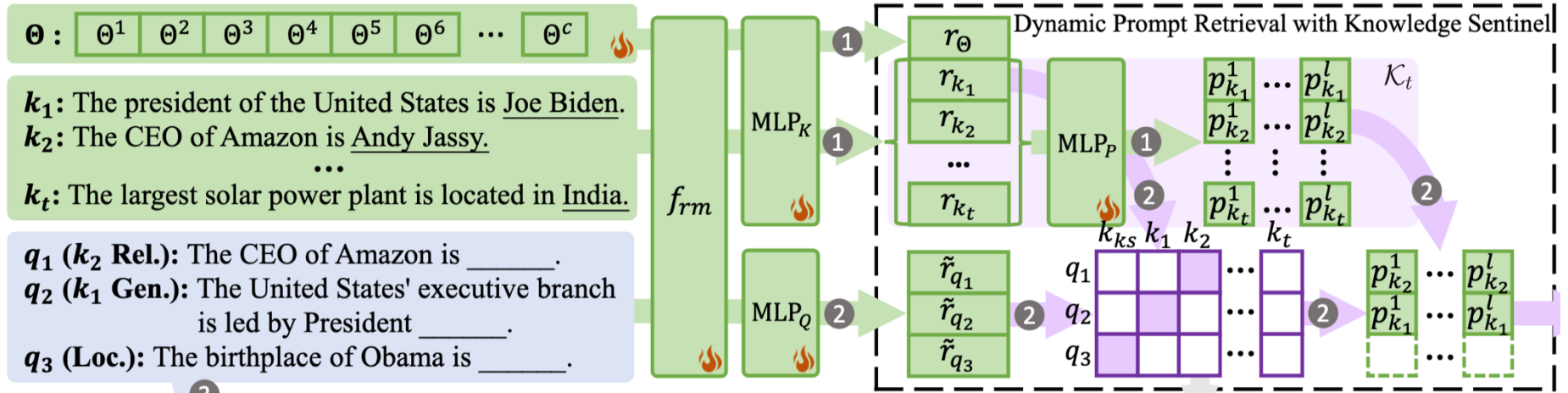
$f_{resp}$  is the reshape operation that maps the vector into a matrix with shape

$\mathcal{K}_t = \mathcal{K}_{t-1} \cup \{(r_{k_t}, p_{k_t})\}$  where  $(r_{k_t}, p_{k_t})$  is key-value pair for knowledge retrieval

# Method

## RECIPE (a RetriEval-augmented Continuous Prompt LEarning)

### 2. Dynamic Prompt Retrieval with Knowledge Sentinel



#### 1) Knowledge Sentinel

(1) Dynamic threshold – fixed threshold X account for the knowledge varies  
 → An intermediary leveraged to dynamically compute similarity threshold

(2)  $KS \Theta \in R$  is a trainable word embedding of  $f_{rm}$  with token length  $c$

(3) Knowledge representation space로 끌고오면..  $r_\Theta = \mathbf{MLP}_K(f_{rm}(\Theta))$

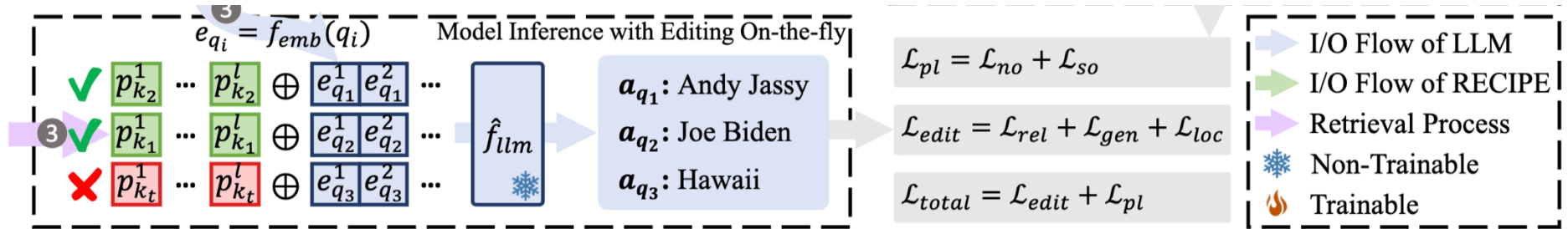
(4) 입력 q에 대한 retrieval process:  $\tilde{r}_q = \mathbf{MLP}_Q(f_{rm}(q))$

$$KS(q) = \begin{cases} p_{k_j} & \tilde{r}_q^T \cdot r_{k_j} > \tilde{r}_q^T \cdot r_\Theta \\ \emptyset & \text{otherwise} \end{cases}$$

# Method

## RECIPE (a RetriEval-augmented Continuous Prompt LEarning)

### 3. Model Inference with Editing On-the-fly



#### 1) Integration Issue

(1) Retrieved continuous prompt + word embedding “q”  $f_{llm} : \mathcal{Q} \mapsto \mathcal{A}$ , where  $\hat{f}_{llm}$  is  $f_{llm}$

$$a_q = \hat{f}_{llm}(p_{k_\tau} \oplus f_{emb}(q))$$

#### (2) Knowledge Editing as a mini-task

- ✗ fine-tuning a specific **prompt encoder** for each mini-task,
- ✓ training **RECIPE modules** that **generate continuous prompts**, ensuring the LLM adheres to the corresponding knowledge.

## RECIPE (a RetriEval-augmented Continuous Prompt LEarning)

### 4. Model Training

**Objective:** Generate continuous prompts and effective retrieval of query-related knowledge for the LLM.

Given a batch of training data consisting of  $b$  editing sample pairs,

$$\{(q_{e_i}, a_{e_i})\}_{i=1}^b \quad \{(q_{g_i}, a_{g_i})\}_{i=1}^b \quad \{(q_{l_i}, a_{l_i})\}_{i=1}^b$$

#### 1) Editing loss

(1) Aim to ensure that the generated continuous prompt guides the LLM to follow the properties  
- **reliability, generality, locality**

$$\mathcal{L}_{rel}^{(i)} = -\log \hat{f}_{llm}(a_{e_i} | p_{k_i} \oplus f_{emb}(q_{e_i})) \quad (10)$$

$$\mathcal{L}_{gen}^{(i)} = -\log \hat{f}_{llm}(a_{g_i} | p_{k_i} \oplus f_{emb}(q_{g_i})) \quad (11)$$

$$\mathcal{L}_{loc}^{(i)} = \text{KL} \left( f_{llm}(q_{l_i}) \parallel \hat{f}_{llm}(p_{k_i} \oplus f_{emb}(q_{l_i})) \right) \quad (12)$$

$$\mathcal{L}_{edit} = \frac{1}{b} \sum_{i=1}^b \left( \mathcal{L}_{rel}^{(i)} + \mathcal{L}_{gen}^{(i)} + \mathcal{L}_{loc}^{(i)} \right).$$

## RECIPE (a RetriEval-augmented Continuous Prompt Learning)

### 4. Model Training

#### 2) Prompt learning loss

(1) Contrastive learning for aligned with the properties of reliability, generality, and locality

For batch of samples,  $\mathcal{L}_{no}^{(i)} = \delta(\tilde{r}_{q_{e_i}}, r_{k_i}, R) + \delta(\tilde{r}_{q_{g_i}}, r_{k_i}, R),$

$$\mathcal{L}_{so}^{(i)} = \delta(\tilde{r}_{q_{l_i}}, r_{\Theta}, R) + \delta(\tilde{r}_{q_{e_i}}, r_{\Theta}, R \setminus k_i) \\ + \delta(\tilde{r}_{q_{g_i}}, r_{\Theta}, R \setminus k_i),$$

$$\mathcal{L}_{pl} = \frac{1}{b} \sum_{i=1}^b (\mathcal{L}_{no}^{(i)} + \mathcal{L}_{so}^{(i)}),$$

where  $R = \{r_{k_i}\}_{i=1}^b \cup \{r_{\Theta}\}$  and  $R \setminus k_i = R \setminus \{r_{k_i}\}$

InfoNCE loss:  $\delta(q, k_+, \{k_i\}_{i=1}^n) = -\log \frac{\exp(q \cdot k_+ / \tau)}{\sum_{i=1}^n \exp(q \cdot k_i / \tau)},$

#### Neighbor-oriented loss + Sentinel-oriented loss



# Experiments

## Editing Performance

# Editing	Type	Editor	ZSRE				CF				RIPE			
			Rel.	Gen.	Loc.	Avg.	Rel.	Gen.	Loc.	Avg.	Rel.	Gen.	Loc.	Avg.
1	MP	FT	47.86	42.57	93.89	61.44 <sub>(±1.00)</sub>	41.37	26.04	52.25	39.89 <sub>(±0.74)</sub>	41.54	33.89	53.27	42.90 <sub>(±0.33)</sub>
		MEND	73.86	70.33	66.10	70.10 <sub>(±0.96)</sub>	81.06	67.15	77.13	75.11 <sub>(±0.62)</sub>	66.37	29.37	29.68	41.81 <sub>(±0.91)</sub>
		ROME	53.49	51.58	93.96	66.34 <sub>(±0.69)</sub>	41.07	21.82	91.85	51.58 <sub>(±0.82)</sub>	48.33	27.08	42.48	39.30 <sub>(±0.89)</sub>
		MEMIT	49.67	49.36	91.87	63.64 <sub>(±0.61)</sub>	45.40	29.25	92.93	55.86 <sub>(±0.39)</sub>	58.37	29.54	38.67	42.19 <sub>(±0.39)</sub>
		MALMEN	46.37	47.75	33.73	42.62 <sub>(±0.43)</sub>	52.45	42.31	36.58	43.78 <sub>(±0.58)</sub>	51.53	33.86	20.45	35.28 <sub>(±1.05)</sub>
		WILKE	50.71	48.52	93.31	64.18 <sub>(±0.55)</sub>	40.07	21.92	91.70	51.23 <sub>(±0.45)</sub>	47.85	27.90	38.50	38.08 <sub>(±1.02)</sub>
	AP	TP	86.35	83.98	86.34	85.56 <sub>(±0.53)</sub>	91.41	68.61	38.94	66.32 <sub>(±1.18)</sub>	76.98	55.10	51.29	61.13 <sub>(±0.48)</sub>
	RB	GRACE	99.20	33.23	99.82	77.42 <sub>(±0.78)</sub>	98.65	11.42	98.73	69.60 <sub>(±0.66)</sub>	98.13	28.45	99.75	75.44 <sub>(±0.65)</sub>
		R-ROME	51.87	49.40	98.82	66.70 <sub>(±1.54)</sub>	39.46	20.76	97.38	52.54 <sub>(±0.86)</sub>	46.15	23.95	92.99	54.37 <sub>(±0.96)</sub>
		LTE	98.97	97.29	85.90	94.05 <sub>(±0.15)</sub>	98.12	97.13	92.20	95.81 <sub>(±1.21)</sub>	98.49	88.09	85.79	90.79 <sub>(±0.61)</sub>
RECIPE		<b>99.40</b>	<b>99.01</b>	<b>99.96</b>	<b>99.46</b> <sub>(±0.07)</sub>	<b>98.78</b>	<b>98.78</b>	<b>99.01</b>	<b>98.86</b> <sub>(±0.39)</sub>	<b>99.36</b>	<b>89.56</b>	<b>99.78</b>	<b>96.24</b> <sub>(±0.95)</sub>	
1000	MP	FT	14.66	12.61	2.69	9.99 <sub>(±1.00)</sub>	6.94	0.68	3.48	3.70 <sub>(±0.09)</sub>	7.91	2.13	1.82	3.95 <sub>(±0.40)</sub>
		MEND	0.04	0.02	0.00	0.02 <sub>(±0.01)</sub>	0.01	0.00	0.02	0.01 <sub>(±0.00)</sub>	0.00	0.02	0.02	0.02 <sub>(±0.00)</sub>
		ROME	1.54	1.48	0.63	1.22 <sub>(±0.90)</sub>	0.15	0.13	0.12	0.14 <sub>(±0.03)</sub>	0.02	0.01	0.03	0.02 <sub>(±0.01)</sub>
		MEMIT	0.18	0.22	0.14	0.18 <sub>(±0.07)</sub>	0.09	0.05	0.99	0.38 <sub>(±0.18)</sub>	0.02	0.02	0.03	0.02 <sub>(±0.01)</sub>
		MALMEN	32.03	28.50	28.14	29.56 <sub>(±1.33)</sub>	15.80	16.41	22.53	18.25 <sub>(±0.22)</sub>	42.33	38.45	38.52	39.77 <sub>(±0.97)</sub>
		WILKE	15.19	12.60	25.31	17.70 <sub>(±1.32)</sub>	13.22	12.28	43.09	22.86 <sub>(±0.64)</sub>	15.19	14.25	10.99	13.48 <sub>(±1.15)</sub>
	AP	TP	44.72	41.38	4.38	30.16 <sub>(±1.04)</sub>	64.70	32.50	11.63	36.28 <sub>(±0.72)</sub>	42.24	26.80	9.87	26.30 <sub>(±1.01)</sub>
	RB	GRACE	42.04	33.42	96.73	57.40 <sub>(±0.68)</sub>	52.75	12.86	91.02	52.21 <sub>(±0.85)</sub>	38.03	30.10	91.24	53.12 <sub>(±0.61)</sub>
		R-ROME	48.73	36.49	94.09	59.77 <sub>(±0.77)</sub>	35.64	14.03	87.94	45.87 <sub>(±0.91)</sub>	41.49	16.96	68.98	42.48 <sub>(±1.21)</sub>
		LTE	93.03	91.14	84.42	89.53 <sub>(±1.16)</sub>	95.87	95.27	89.35	93.50 <sub>(±0.26)</sub>	94.53	84.52	80.44	86.50 <sub>(±0.75)</sub>
RECIPE		<b>96.30</b>	<b>95.27</b>	<b>99.98</b>	<b>97.18</b> <sub>(±0.50)</sub>	<b>96.37</b>	<b>96.04</b>	<b>93.66</b>	<b>95.35</b> <sub>(±0.61)</sub>	<b>95.60</b>	<b>85.53</b>	<b>92.35</b>	<b>91.16</b> <sub>(±1.28)</sub>	

## Experiments

### Linguistic Capability

Editor	CSQA	MMLU	ANLI	SQUAD-2	Average
N/A	38.91	41.54	34.04	36.43	37.73
FT	19.27	29.93	33.33	0.59	20.78
MEND	20.31	24.68	33.07	0.04	19.52
ROME	19.97	23.03	33.47	0.41	19.22
MEMIT	19.68	23.23	33.39	0.01	19.08
TP	19.62	22.84	33.37	0.96	19.20
GRACE	38.60	41.20	33.93	36.28	37.50
R-ROME	38.50	41.12	33.90	36.31	37.46
MALMEN	20.85	24.83	33.03	0.27	19.75
LTE	19.45	23.21	33.41	25.25	25.33
WILKE	19.87	23.37	33.37	0.07	19.17
RECIPE	<b>38.76</b>	<b>41.40</b>	<b>34.13</b>	<b>36.50</b>	<b>37.70</b>

Table 2: Performance of LLAMA-2 after 1,000 edits. “N/A” denotes performance without any edits. Bold font highlights the optimal post-editing performance.

# Experiments

## Further..

Type	Editor	Edit Time	Infer. Time	Total Time
MP	FT	1.7205	0.0589	1.7794
	MEND	0.0987	0.0590	0.1577
	ROME	17.1639	0.0586	17.2225
	MEMIT	33.6631	0.0591	33.7222
	MALMEN	2.3418	0.0589	2.4007
	WILKE	38.7165	0.0587	38.7752
AP	TP	5.9061	0.0615	5.9676
RB	GRACE	12.5343	0.0936	12.6279
	R-ROME	17.3135	0.0606	17.3741
	LTE	0.0076	0.0634	0.0710
	RECIPE	0.0078	0.0598	0.0676

Table 3: Average time (s) taken for a single edit and model inference after 10,000 edits.

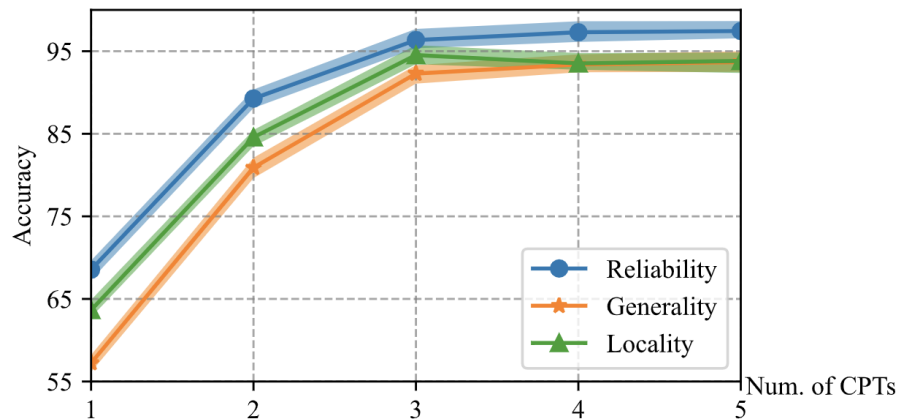


Figure 3: Impact of the number of CPTs on editing performance of RECIPE.

Settings	100 Edits			1000 Edits		
	Rel.	Gen.	Loc.	Rel.	Gen.	Loc.
N/A	27.30	26.07	100.00	27.30	26.07	100.00
RECIPE	<b>97.29</b>	<b>93.74</b>	97.38	<b>96.05</b>	<b>92.34</b>	95.36
- CPT	27.42	26.18	<b>99.98</b>	27.38	26.15	<b>99.97</b>
- KS	95.55	89.10	92.45	94.01	86.63	88.55
- BOTH	27.41	26.17	99.96	27.35	26.12	99.94

Table 4: Ablation study of RECIPE.

## Conclusion

---

- 1) RAG의 강력함 (뭔가.. Converge 하는 것 같기도 하고, 평균 이상의 성능을 보이기도 함)
- 2) Knowledge Editing은 여러 도메인의 기술이 접목되는 영역
- 3) 과거의 아이디어도 다시 빛을 발휘할 수 있음
- 4) LTE & RECIPE의 약점은 무엇일까?

## Q&A