

# Knowledge Editing

Large Language Models: Introduction and Recent Advances

ELL881 · AIL821



Tanmoy Chakraborty  
Associate Professor, IIT Delhi  
<https://tanmoychak.com/>



# Mistral AI announces Pixtral !

Announced on  
September 17,  
2024

[Mistral AI Blog](#)

Pixtral 12B is the first multimodal model by Mistral AI.

Pixtral consists of a **new 400M parameter vision encoder trained from scratch**, and a **12B parameter multimodal decoder based on Mistral Nemo**. Pixtral supports multiple images in the long context window of 128k tokens.

# Pixtral 12B



Pixtral supports variable image sizes and aspect ratios. It understands both natural images and documents, **achieving 52.5% on the MMMU reasoning benchmark, surpassing a number of larger models**

# LLMs Need Regular Updates



<https://arxiv.org/pdf/2310.16218>



# Issues with Finetuning

- **Computationally expensive**

LLaMA was trained for **21 days** on **2,048 A100 GPUs**, costing over **\$2.4M** and emitting over **1,000 tons** of CO2 [Hartvigsen et al., NeurIPS'23]

- **Unconstrained editing**

Fine-tuning LLMs alters the pre-trained parameters without constraints, leading to the overfitting problem [Wortsman et al, CVPR'22]



# Knowledge Editing: Definition

**Knowledge Edit:** A factual information can be presented as a triplet  $\langle s, r, o \rangle$

where,

s : subject      r : relation      o : object

**Example: “The Space Needle is located in the city of Seattle”**

s : subject :    The Space Needle

r : relation :    Location

o : object :     Seattle

It can be expressed as  $\langle$ The space Needle, Location, Seattle $\rangle$



# Goal

- Update the objective from 'o' to 'o\*' i.e. < s, r, o > to < s, r, o\* >

Example:

< The space Needle, Location, Seattle > -> <The space Needle, Location, Goa >

< The US, President, Barack Obama > -> < The US, President, Joe Biden >

< Earth, Highest point, Mount Everest > -> < Earth, Highest point, Mount K2 >

< Adult human, No. of bones, 206 > -> < Adult human, No. of bones, 300 >

< Water, Boiling point, 100 °C > -> < Water, Boiling point, 200 °C >



# Conditions for Successful Edit

**Condition 1: Reliability** : It expects  $f_{\text{edit}}(\theta')$  to select  $y_{\text{alt}}$  instead of  $y_{\text{old}}$  for a given input triplet  $t$  from  $D_x$  where  $D_x$  is the set of targeted queries

Example:

For  $t$  : < The US, President, **Barack Obama** > to < The US, President, **Joe Biden** >

Input query : The president of the US is ...

Expected output : **Barack Obama** **Joe Biden**



# Conditions for Successful Edit

**Condition 2: Generalization:** It expects  $f_{\text{edit}}(\theta')$  to select  $y_{\text{alt}}$  instead of  $y_{\text{old}}$  for the paraphrased versions of inputs from  $D_x$ , denoted by  $P_x$

Example:

For  $t$ : < The US, President, **Barack Obama** > to < The US, President, **Joe Biden** >

Input query : The president of the US is

Expected output : **Barack Obama** **Joe Biden**

Input query : Who is the president of the US?

Expected output : **Barack Obama** **Joe Biden**





# Conditions for Successful Edit

**Condition 3: Localization:** It expects  $f_{\text{edit}}(\theta')$  to select  $y_{\text{old}}$  instead of  $y_{\text{alt}}$  for a given input triplet  $t$  from  $O_x$  where  $O_x$  is the set of non-targeted queries.

## Example:

For  $t$  : < Russia, President, Vladimir Putin >

Input query : The president of Russia is

Expected output : Vladimir Putin

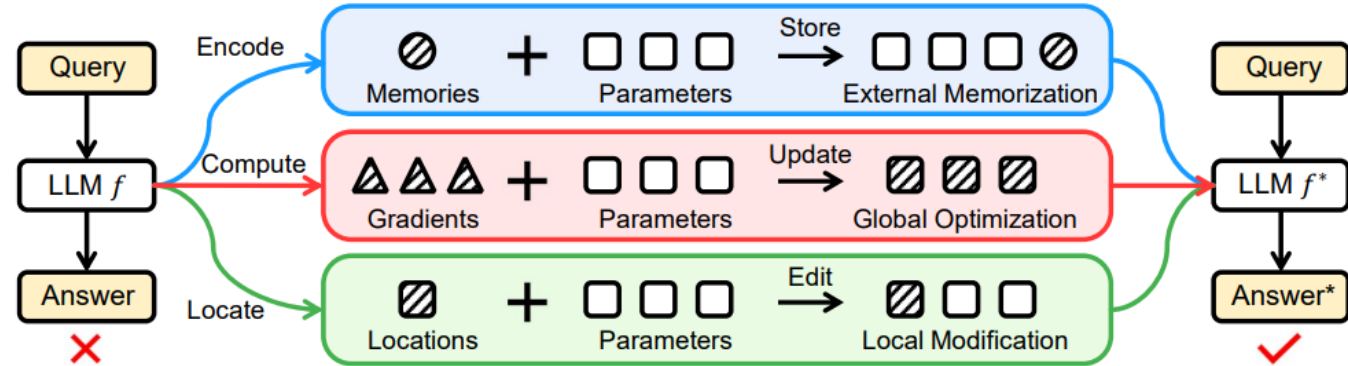
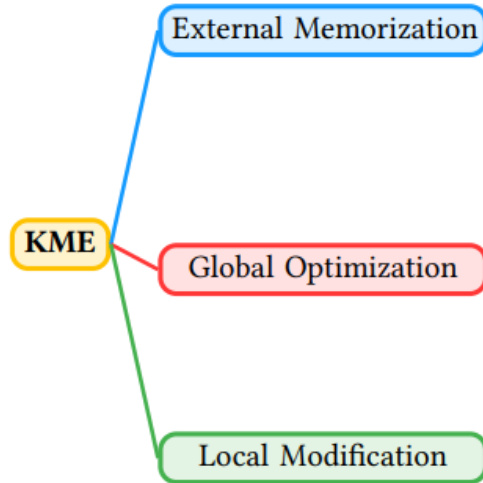
Input query : Who is the president of Russia?

Expected output : Vladimir Putin



# Taxonomy

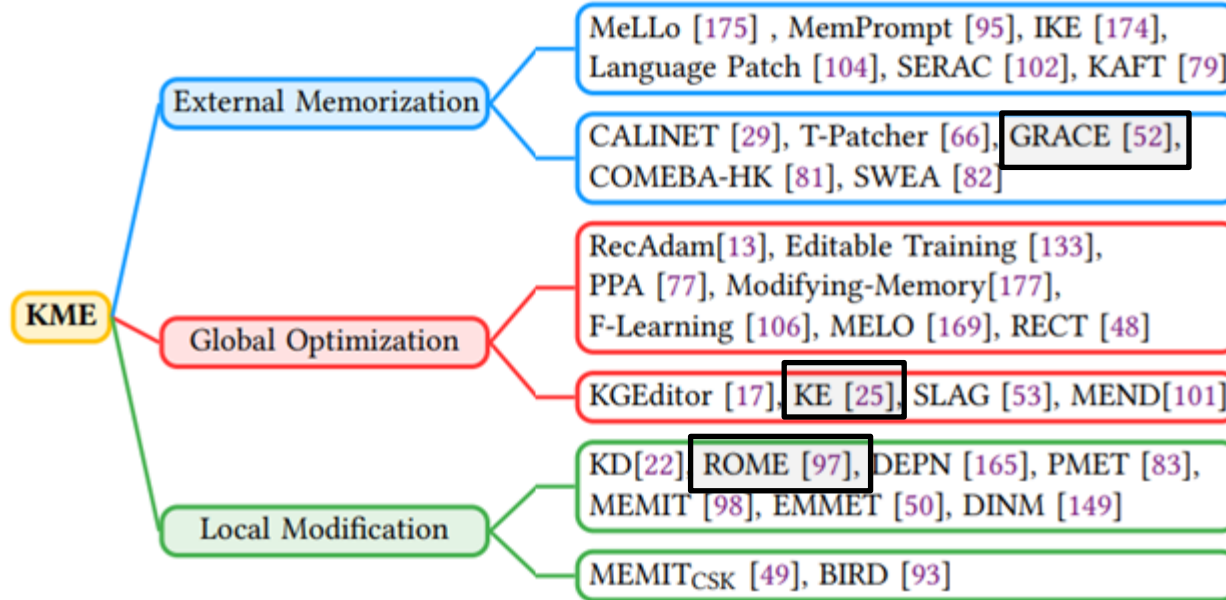
**Where?** (e.g., external parameters or internal weights)  
**How?** (e.g., via optimization or direct incorporation)



<https://arxiv.org/pdf/2310.16218>



# Taxonomy



<https://arxiv.org/pdf/2310.16218>



# KNOWLEDGE EDITOR

## Editing Factual Knowledge in Language Models

Nicola De Cao <sup>1,2</sup>, Wilker Aziz <sup>1</sup>, Ivan Titov <sup>1,2</sup>

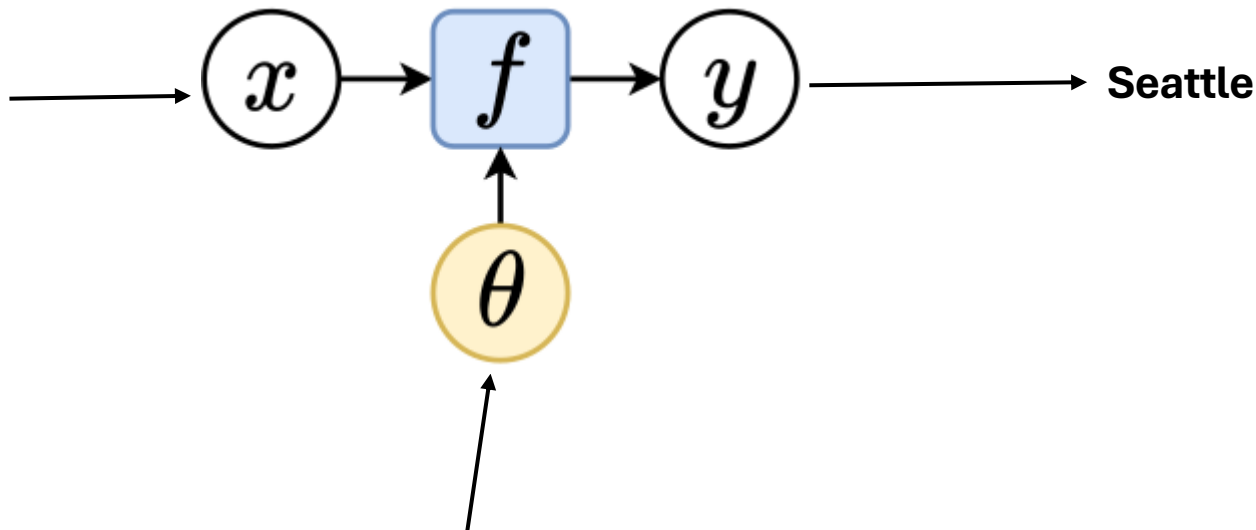
<sup>1</sup>University of Amsterdam, <sup>2</sup>University of Edinburgh

{ nicola.decao, w.aziz, titov } @uva.nl



# Method: KNOWLEDGE EDITOR

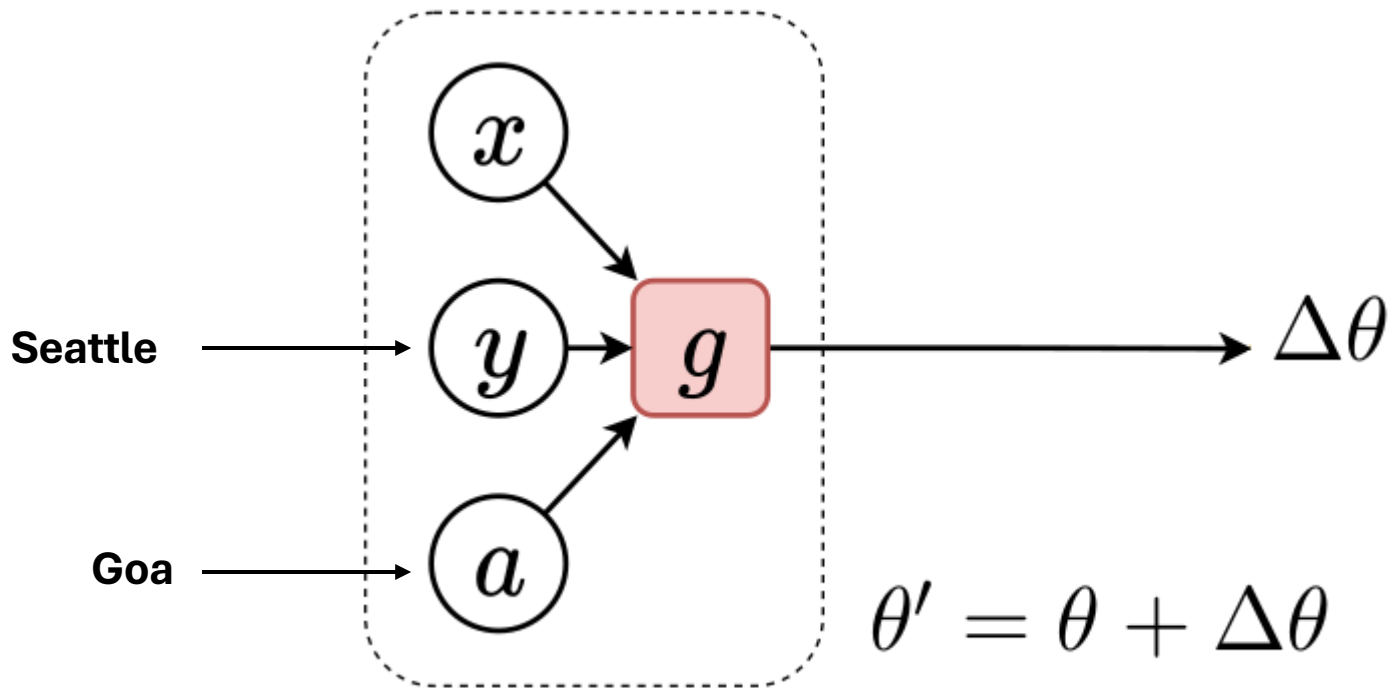
The Space Needle is located in the city of



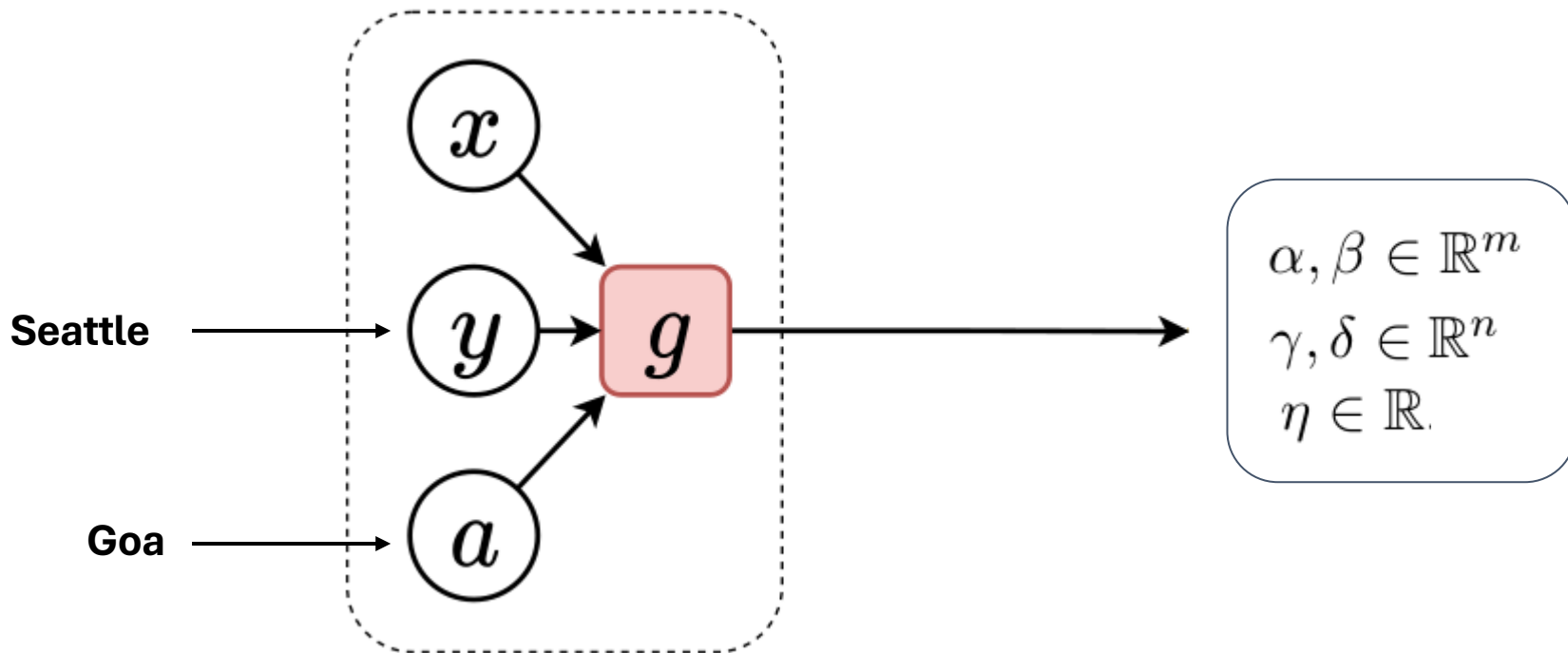
The parameters of a **pre-trained** model



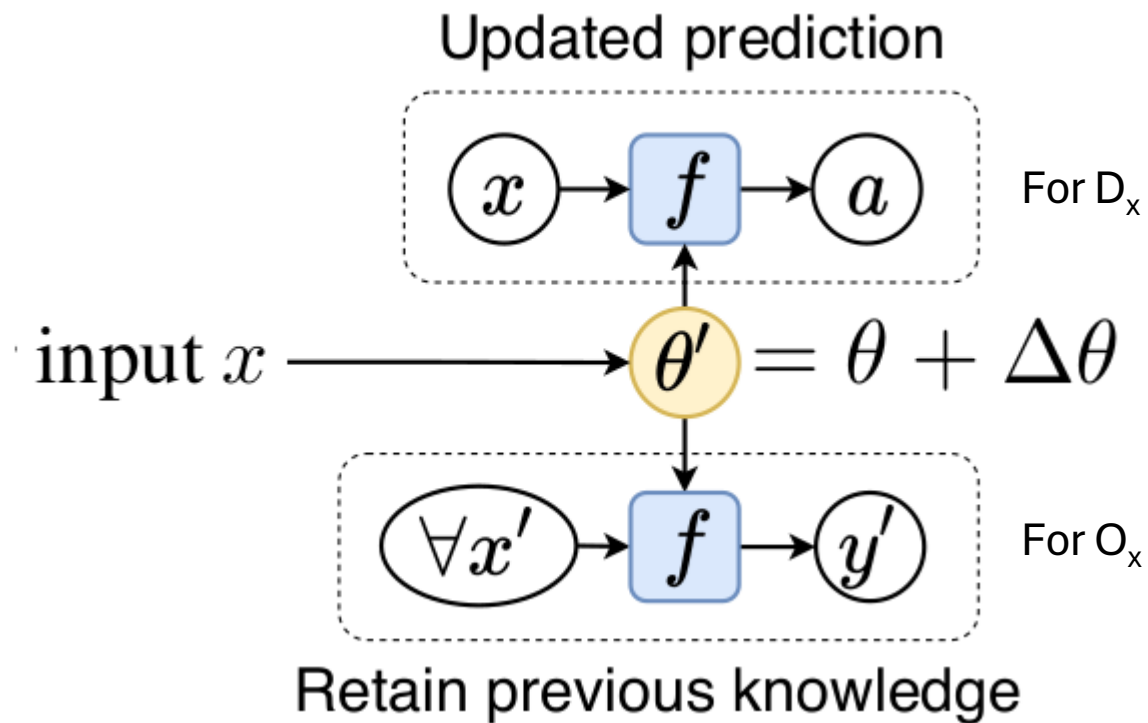
# Method: KNOWLEDGE EDITOR



# Method: KNOWLEDGE EDITOR



# Method: KNOWLEDGE EDITOR





# Training Hyper-network

$$\min_{\phi} \sum_{\hat{x} \in \mathcal{P}^x} \mathcal{L}(\theta'; \hat{x}, a)$$

$$\text{s.t. } \mathcal{C}(\theta, \theta', f; \mathcal{O}^x) < m$$

$$\mathcal{C}_{KL}(\theta, \theta', f; \mathcal{O}^x) = \sum_{x' \in \mathcal{O}^x} \sum_{c \in \mathcal{Y}} p_{Y|X}(c|x', \theta) \log \frac{p_{Y|X}(c|x', \theta)}{p_{Y|X}(c|x', \theta')}$$

$$\mathcal{C}_{L_p}(\theta, \theta', f; \mathcal{O}^x) = (\sum_i |\theta_i - \theta'_i|^p)^{1/p}$$



# Training Hyper-network

$$\min_{\phi} \sum_{\hat{x} \in \mathcal{P}^x} \mathcal{L}(\theta'; \hat{x}, a)$$

$$\text{s.t. } \mathcal{C}(\theta, \theta', f; \mathcal{O}^x) < m$$

$$\mathcal{C}_{KL}(\theta, \theta', f; \mathcal{O}^x) = \sum_{x' \in \mathcal{O}^x} \sum_{c \in \mathcal{Y}} p_{Y|X}(c|x', \theta) \log \frac{p_{Y|X}(c|x', \theta)}{p_{Y|X}(c|x', \theta')}$$

$$\mathcal{C}_{L_p}(\theta, \theta', f; \mathcal{O}^x) = (\sum_i |\theta_i - \theta'_i|^p)^{1/p}$$

$$\Delta W = \sigma(\eta) \cdot (\hat{\alpha} \odot \nabla_W \mathcal{L}(W; x, a) + \hat{\beta}) ,$$

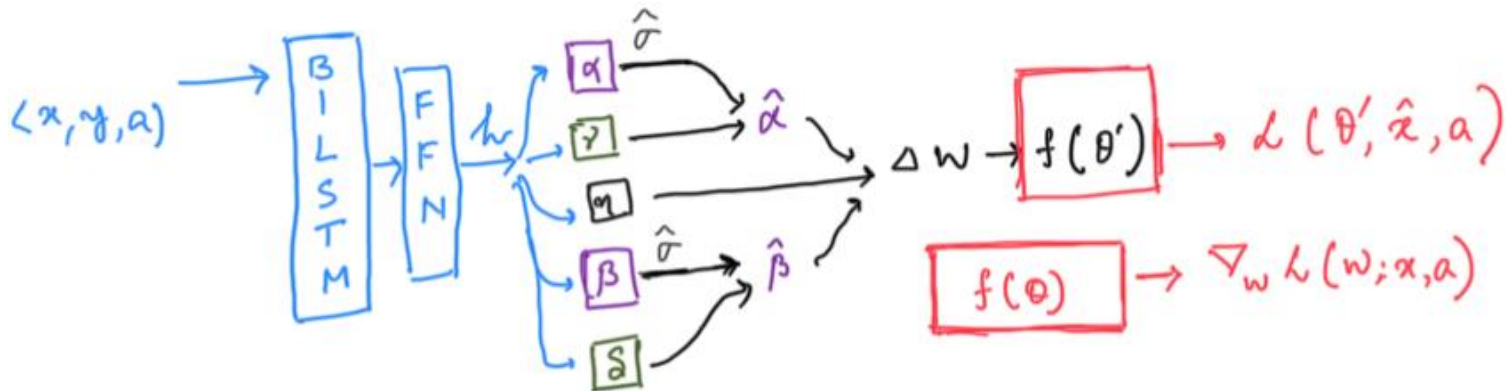
$$\text{with } \hat{\alpha} = \hat{\sigma}(\alpha) \gamma^\top \text{ and } \hat{\beta} = \hat{\sigma}(\beta) \delta^\top ,$$

$\hat{\sigma}$  : Softmax

$\sigma$  : Sigmoid



# Training Hyper-network



$$\min_{\phi} \sum_{\hat{x} \in \mathcal{P}^x} \mathcal{L}(\theta'; \hat{x}, a)$$

$$\text{s.t. } \mathcal{C}(\theta, \theta', f; \mathcal{O}^x) < m$$

$$\mathcal{C}_{KL}(\theta, \theta', f; \mathcal{O}^x) = \sum_{x' \in \mathcal{O}^x} \sum_{c \in \mathcal{Y}} p_{Y|X}(c|x', \theta) \log \frac{p_{Y|X}(c|x', \theta)}{p_{Y|X}(c|x', \theta')}$$

$$\mathcal{C}_{L_p}(\theta, \theta', f; \mathcal{O}^x) = (\sum_i |\theta_i - \theta'_i|^p)^{1/p}$$

$$\Delta W = \sigma(\eta) \cdot (\hat{\alpha} \odot \nabla_W \mathcal{L}(W; x, a) + \hat{\beta}),$$

$$\text{with } \hat{\alpha} = \hat{\sigma}(\alpha) \gamma^T \text{ and } \hat{\beta} = \hat{\sigma}(\beta) \delta^T,$$

$\hat{\sigma}$  : Softmax

$\sigma$  : Sigmoid



# Evaluation: KNOWLEDGE EDITOR

BERT-base on FEVER dataset

Method	Fact-Checking		
	Success rate $\uparrow$	Retain acc $\uparrow$	Equiv. acc $\uparrow$
Fine-tune (1st layer)	100.0	99.44	42.24
Fine-tune (all layers)	100.0	86.95	95.58
Zhu et al. (1st layer)	100.0	99.44	40.30
Zhu et al. (all layers)	100.0	94.07	83.30
Ours $\mathcal{C}_{L_2}$	99.10	45.10	99.01
KNOWLEDGEEDITOR	98.80	98.14	82.69
+ loop <sup>†</sup>	100.0	97.78	81.57
+ $\mathcal{P}^x$ <sup>‡</sup>	98.50	98.55	95.25
+ $\mathcal{P}^x$ + loop <sup>‡</sup>	100.0	98.46	94.65

*success rate*: how much  $g$  successfully updates the knowledge in  $\theta'$ , measured as accuracy of revised predictions for inputs in  $\mathcal{D}$ ;

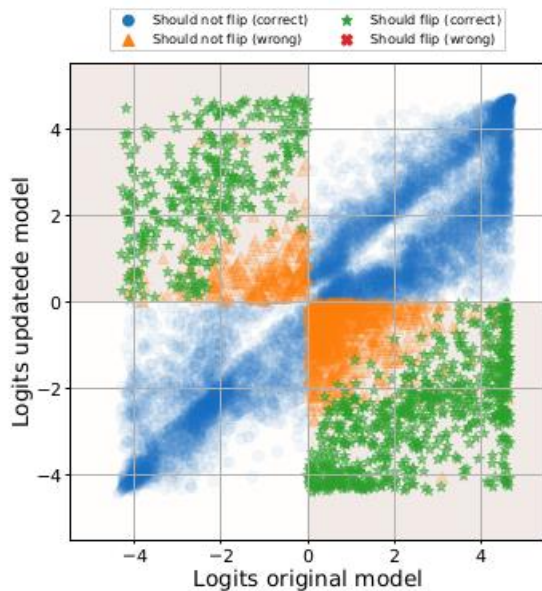
*retain accuracy*: how well  $\theta'$  retains the original predictions of  $f$ , measured as accuracy wrt input-output pairs in sets  $\mathcal{O}^x$ ;

*equivalence accuracy*: how consistent the predictions of the revised model  $\theta'$  are for semantically equivalent inputs, measured as accuracy of the revised predictions for all  $\mathcal{P}^x$ ;



# Evaluation: KNOWLEDGE EDITOR

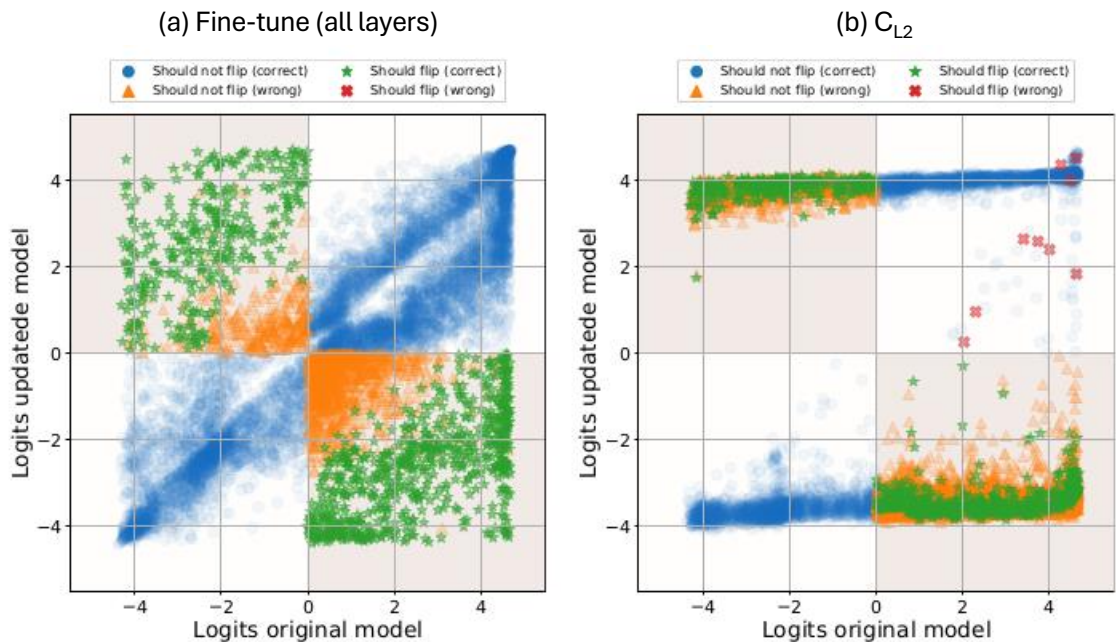
(a) Fine-tune (all layers)



Distribution of logits in the original model and the updated model on FEVER dataset. (a) Fine-tune all layers causes ample amount of errors. (b)  $C_{L_2}$  is able to alter the predictions but fails at retaining the updates. (c) KNOWLEDGE - EDITOR with  $C_{KL}, P_x$  has minimal errors and is able to achieve correct output with a small perturbation.



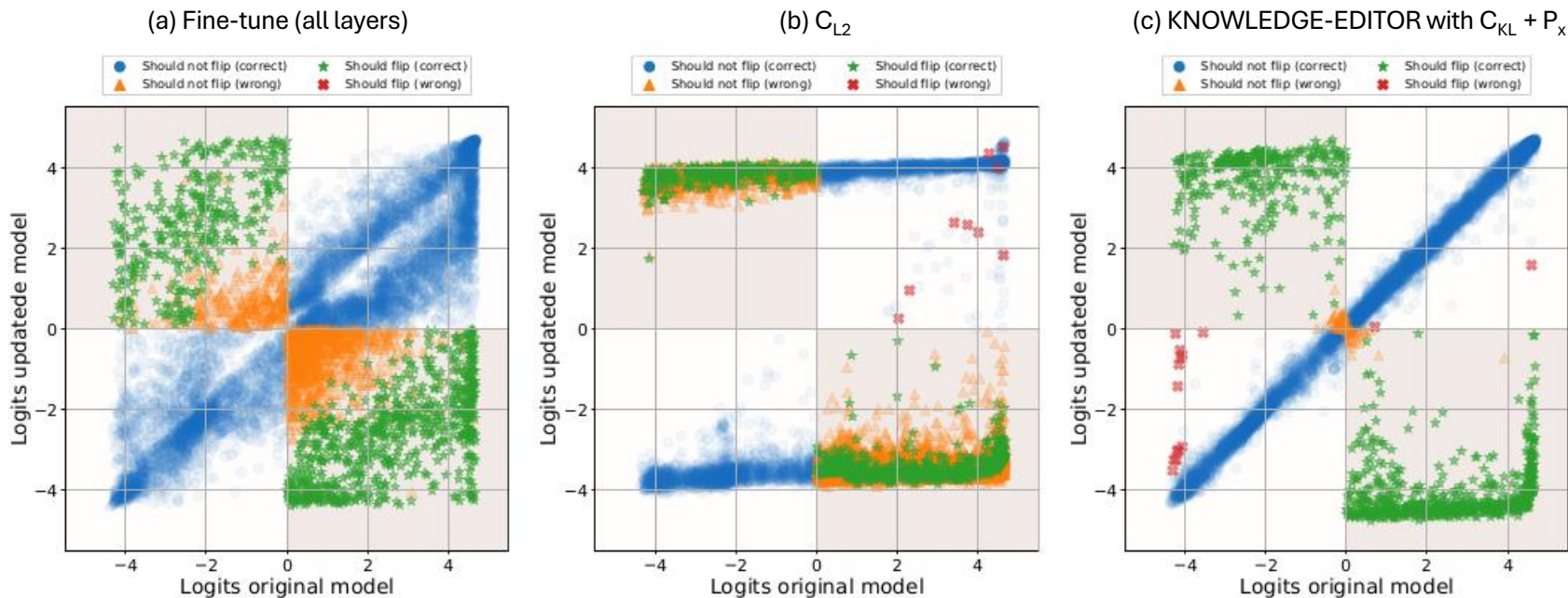
# Evaluation: KNOWLEDGE EDITOR



Distribution of logits in the original model and the updated model on FEVER dataset. (a) Fine-tune all layers causes ample amount of errors. (b)  $C_{L_2}$  is able to alter the predictions but fails at retaining the updates. (c) KNOWLEDGE EDITOR with  $C_{KL}, P_x$  has minimal errors and is able to achieve correct output with a small perturbation.



# Evaluation: KNOWLEDGE EDITOR



Distribution of logits in the original model and the updated model on FEVER dataset. (a) Fine-tune all layers causes ample amount of errors. (b)  $C_{L_2}$  is able to alter the predictions but fails at retaining the updates. (c) KNOWLEDGE-EDITOR with  $C_{KL}, P_x$  has minimal errors and is able to achieve correct output with a small perturbation.



# GRACE

---

## Aging with GRACE: Lifelong Model Editing with Discrete Key-Value Adaptors

---

**Thomas Hartvigsen**  
University of Virginia, MIT  
hartvigsen@virginia.edu

**Swami Sankaranarayanan**  
Sony AI  
swami.sankaranarayanan@sony.com

**Hamid Palangi**  
Microsoft Research  
hpalangi@microsoft.com

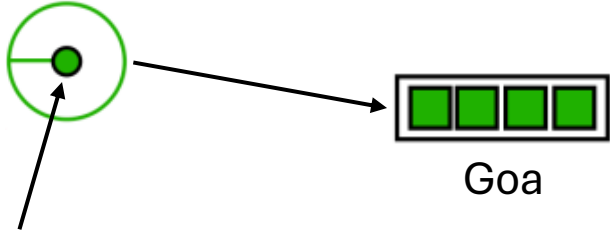
**Yoon Kim**  
MIT  
yoonkim@mit.edu

**Marzyeh Ghassemi**  
MIT  
mghassem@mit.edu

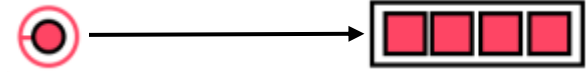




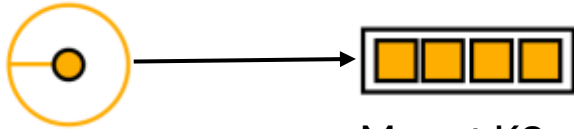
# Method: GRACE



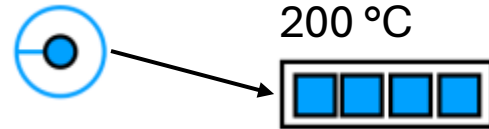
Where is the Space Needle located?



How many bones are there in an adult human?



What is the highest point on Earth?



What is the boiling point of water?



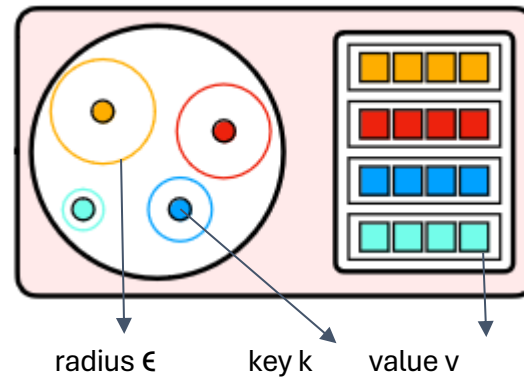
# GRACE Adaptor

- Uses a cache-like storage for a layer, called **GRACE Adaptor**
- An Adaptor contains **two** components:
  - (1) a codebook **C**
  - (2) a deferral mechanism



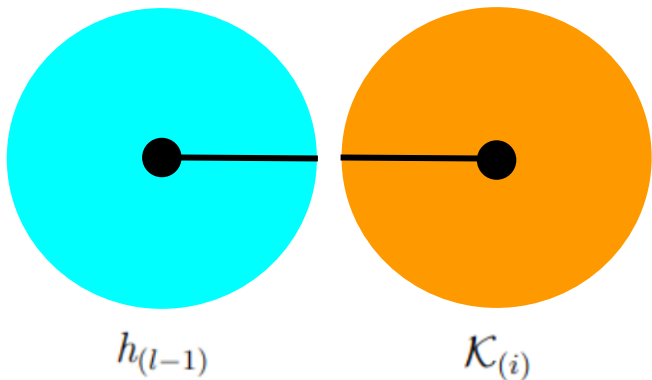
# GRACE: Codebook

- The codebook  $\mathbf{C}$  has three components
  - A set of keys ( $K$ )
  - A set of values ( $V$ )
  - Deferral radii ( $\xi$ )



# GRACE: Deferral Mechanism

- In case of **no match** or **empty** codebook **C**
- A new entry will be created.



$$h_{(l)} = \begin{cases} \text{GRACE}(h_{(l-1)}), & \text{if } \min_i (d(h_{(l-1)}, \mathcal{K}_{(i)})) < \xi_{(i)} \\ f_l(h_{(l-1)}), & \text{otherwise} \end{cases}$$

**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  $\mathcal{C}$

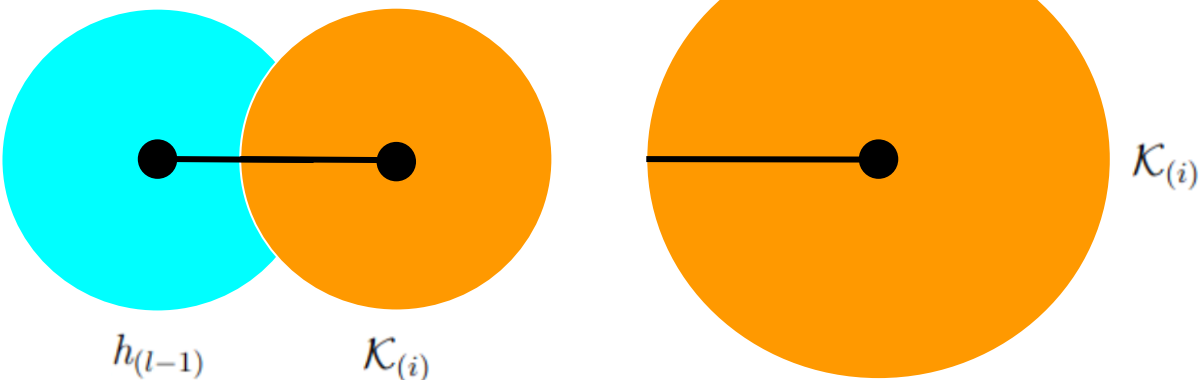
$v_{\text{new}} = \text{finetune on } P_f(y|v_{\text{init}})$

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



# GRACE: Deferral Mechanism

- In case of **collision**, expand if the  $o^*$  is **same**.



$$h^{(l)} = \begin{cases} \text{GRACE}(h^{(l-1)}), & \text{if } \min_i (d(h^{(l-1)}, \mathcal{K}_{(i)})) < \xi(i) \\ f_i(h^{(l-1)}), & \text{otherwise} \end{cases}$$

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

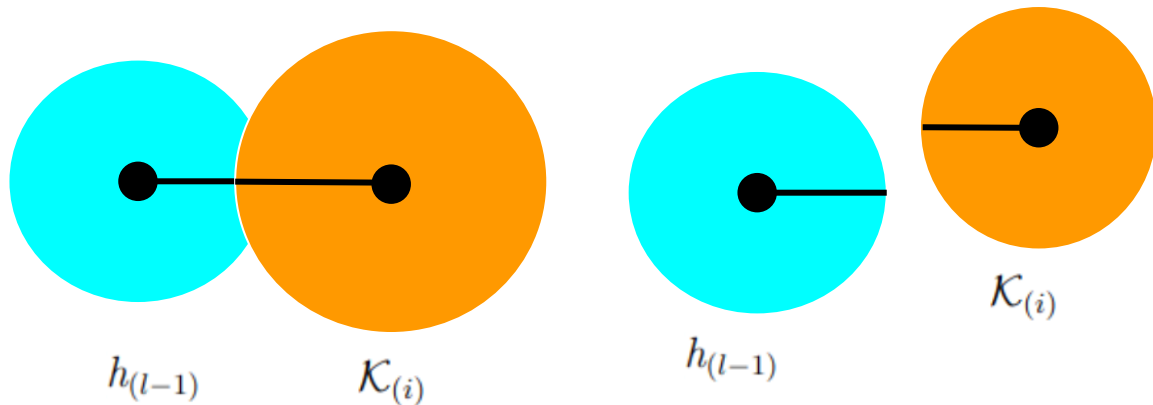
# *Same label*  $\rightarrow$  *Expand*

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



# GRACE: Deferral Mechanism

- In case of **collision**, reduce if the  $o^*$  is **different**.



$$h_{(l)} = \begin{cases} \text{GRACE}(h_{(l-1)}), & \text{if } \min_i (d(h_{(l-1)}, \mathcal{K}_{(i)})) < \xi_{(i)} \\ f_i(h_{(l-1)}), & \text{otherwise} \end{cases}$$

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



# GRACE: Evaluation

Method	zsRE (T5; F1 ↑)				SCOTUS (BERT; Acc ↑)			
	TRR	ERR	Avg.	#E	TRR	ERR	Avg.	#E
FT <b>25</b>	.56	.82	.69	1000	.52	.52	.52	415
FT+EWC <b>19</b>	.51	.82	.66	1000	.67	.50	.58	408
FT+Retrain <b>36</b>	.27	.99	.63	1000	.67	<b>.83</b>	.75	403
MEND <b>30</b>	.25	.27	.26	1000	.19	.27	.23	672
Defer <b>31</b>	<b>.72</b>	.31	.52	1000	.33	.41	.37	506
ROME <b>28</b>	—	—	—	—	—	—	—	—
Memory	.25	.27	.26	1000	.21	.20	.21	780
GRACE	<b>.69</b>	<b>.96</b>	<b>.82</b>	1000	<b>.81</b>	<b>.82</b>	<b>.82</b>	381

#E: the number of edits

#TRR: it evaluates the knowledge retention rate of the edited model.

#ERR: it evaluates the retention rate of old-edits.

- T5-60M, task: context-free QA, zsRE dataset.
- BERT-110M, task: label shifts in legal documents in US, SCOTUS dataset.



# ROME

---

## Locating and Editing Factual Associations in GPT

---

**Kevin Meng\***  
MIT CSAIL

**David Bau\***  
Northeastern University

**Alex Andonian**  
MIT CSAIL

**Yonatan Belinkov†**  
Technion – IIT



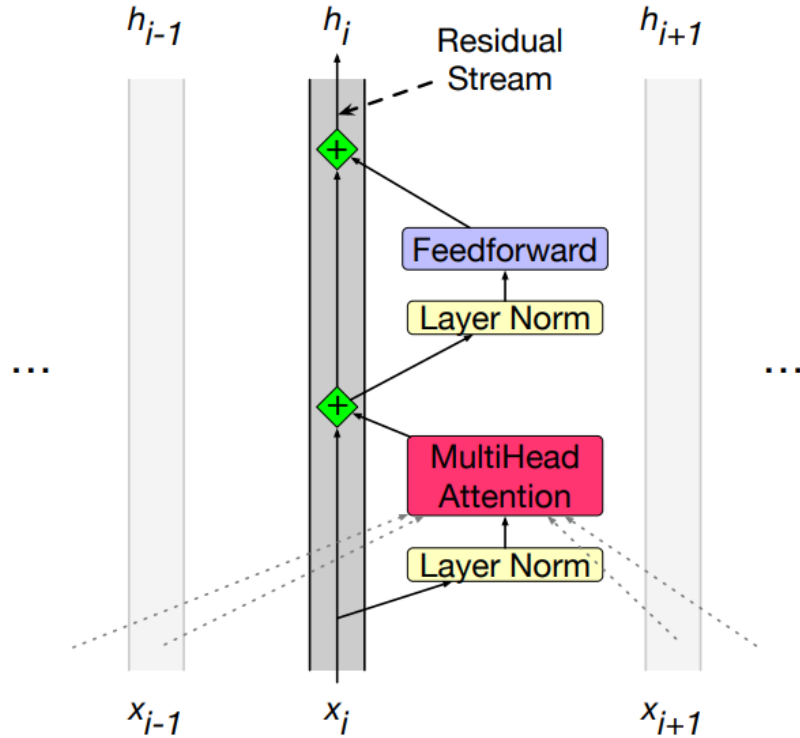


# Method: ROME

- An in-depth analysis on the storage and recall of factual associations in autoregressive transformer language models.
- It provides evidence that factual associations correspond to localized, directly-editable computations.
- It proposes a mechanism to update specific factual associations in feed-forward weights using **ROME**.
- First to find that the mid-layer feedforward modules play a crucial role in storing factual associations and it can be altered with a direct manipulation.



# Residual Stream Perspective of Transformer

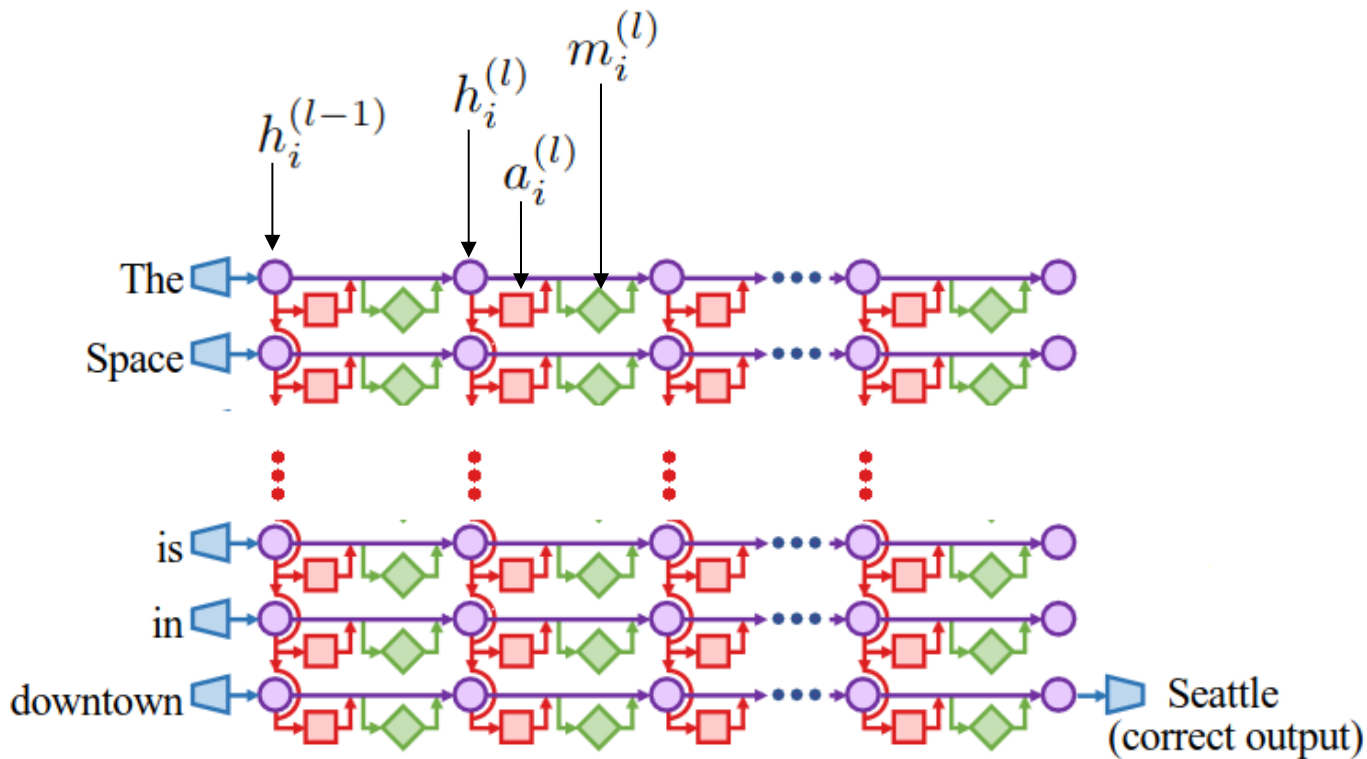


<https://web.stanford.edu/~jurafsky/slp3/9.pdf>



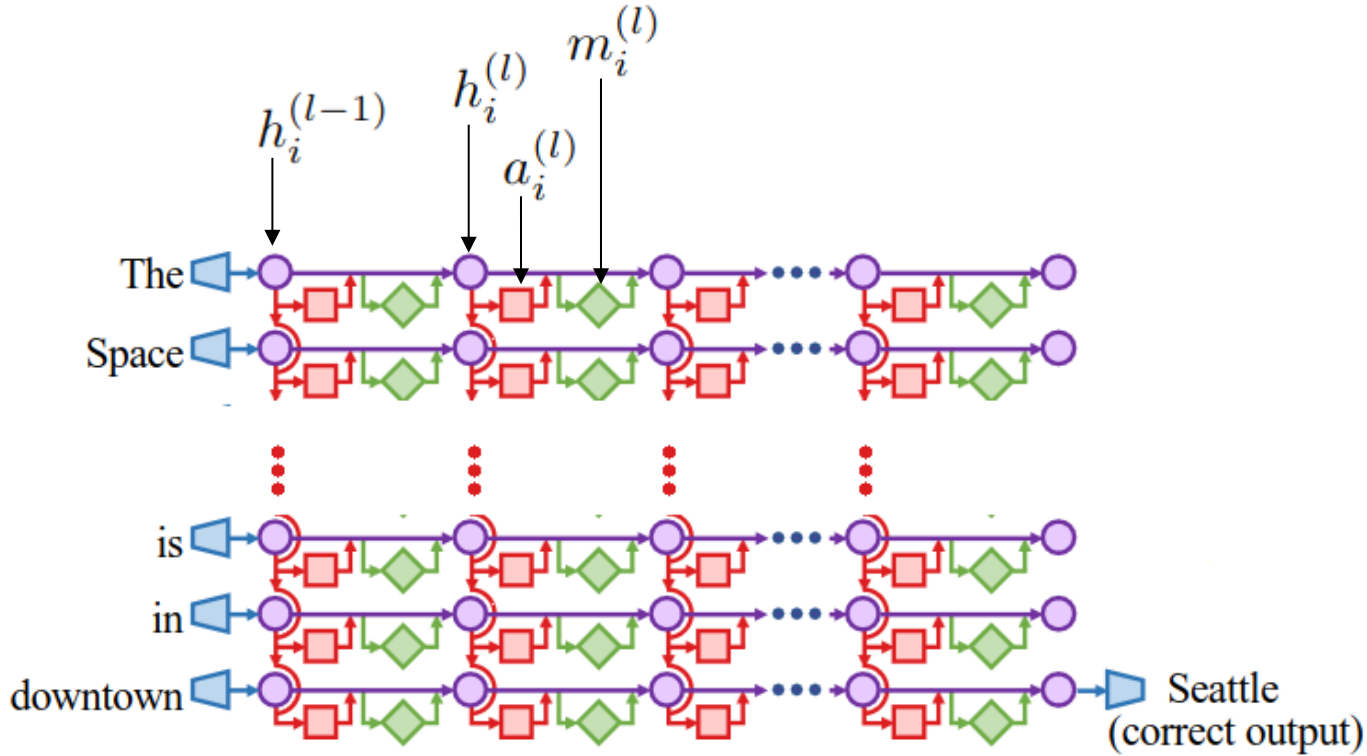
# Method: ROME

- $h_i^{(l)}$  state
- attention  $a_i^{(l)}$
- MLP



# Method: ROME

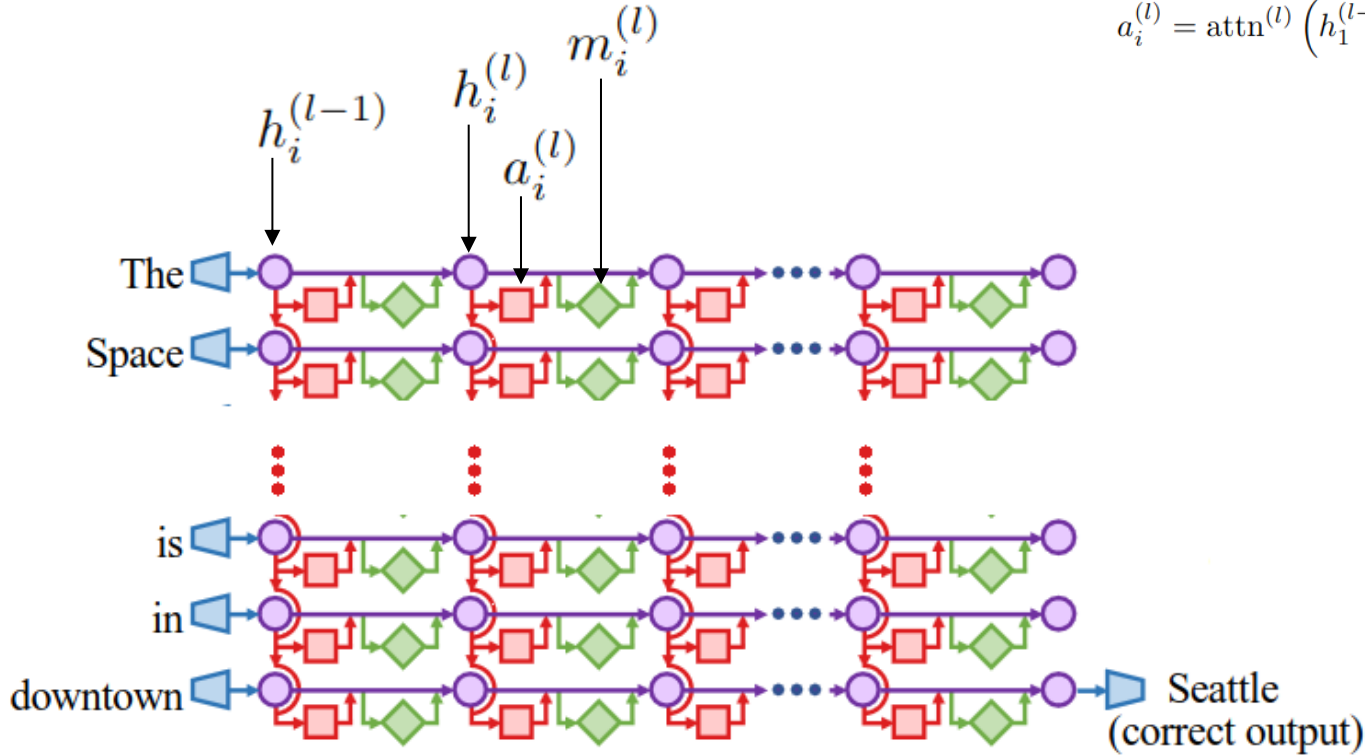
$$h_i^{(l)} = h_i^{(l-1)} + a_i^{(l)} + m_i^{(l)}$$



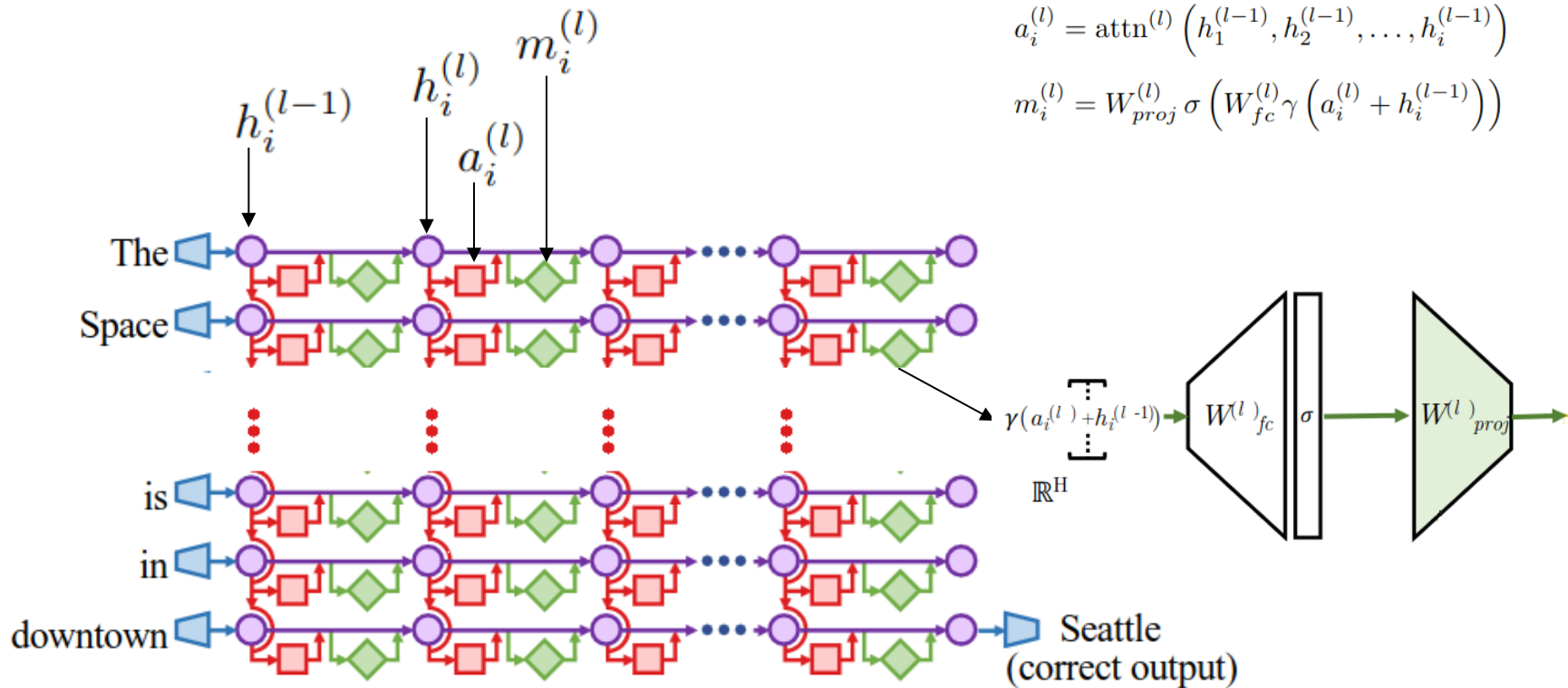
# Method: ROME

$$h_i^{(l)} = h_i^{(l-1)} + a_i^{(l)} + m_i^{(l)}$$

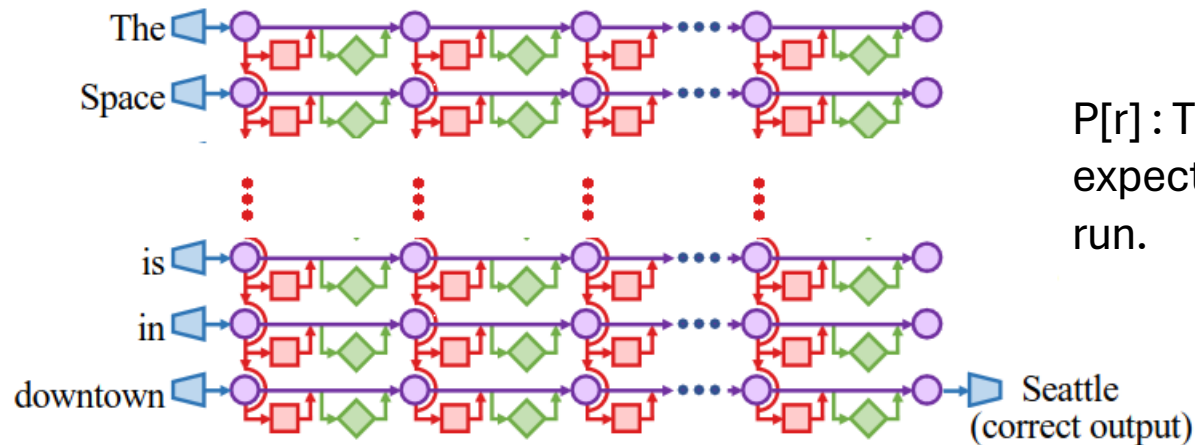
$$a_i^{(l)} = \text{attn}^{(l)}(h_1^{(l-1)}, h_2^{(l-1)}, \dots, h_i^{(l-1)})$$



# Method: ROME



# Method: ROME

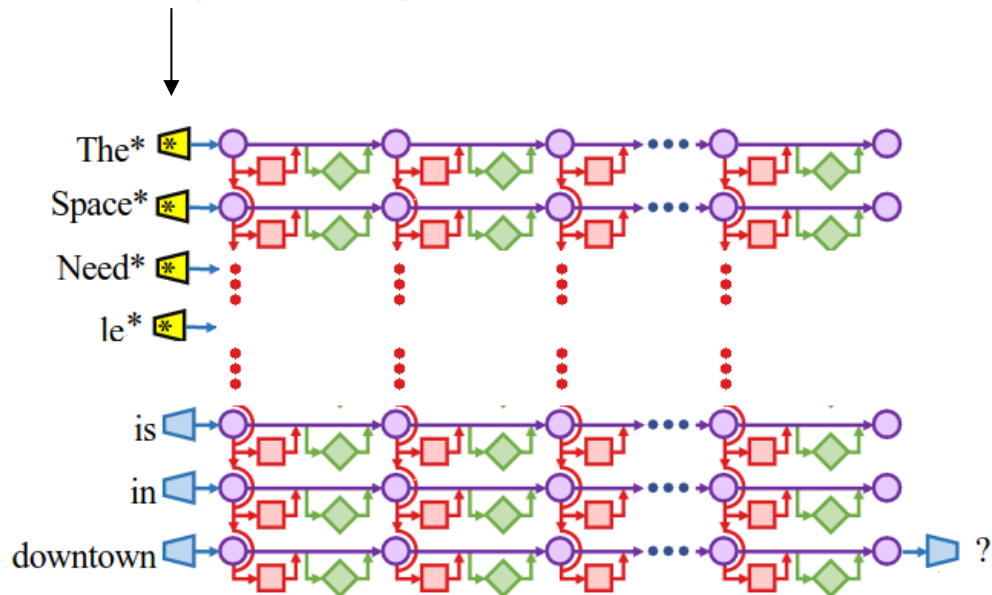


$P[r]$  : The probability of emitting the expected response  $r$  under the **clean** run.



## Method: ROME

$$h_i'^{(0)} = h_i^{(0)} + \epsilon$$



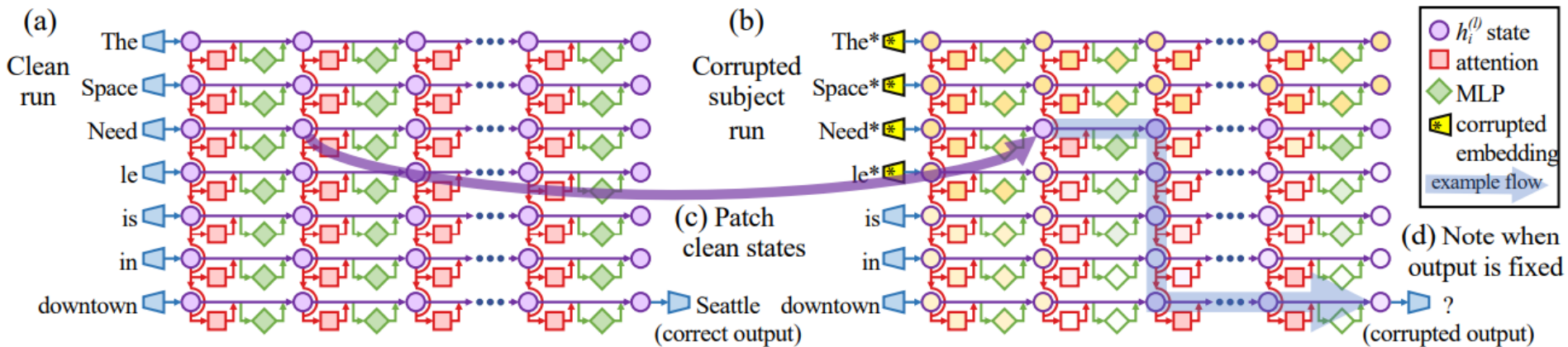
$P_*[r]$ : The probability of emitting the expected response  $r$  under the **corrupted** run.





# Method: ROME

## A corrupted-with-restoration run



# Method: ROME

- The relative importance of a state  $s$  is approximated through the following process
  - **Total Effect (TE)**: difference between the probabilities of generating the true response under clean run and corrupted run

$$TE = P[r] - P_*[r]$$

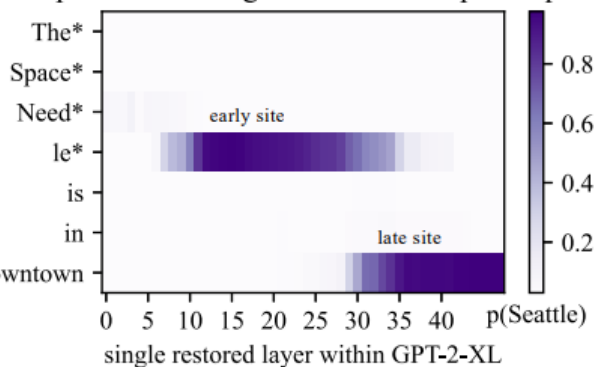
- **Indirect Effect (IE)**:

$$IE = P_*^{(k)}_i[r] - P_*[r]$$

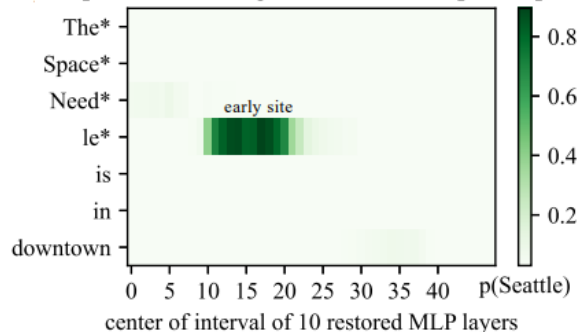


# Effect of Causal Tracing

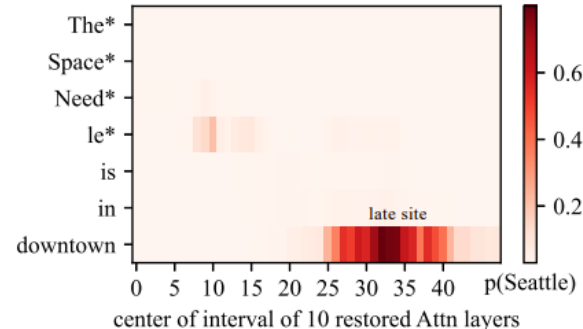
Impact of restoring state after corrupted input



Impact of restoring MLP after corrupted input



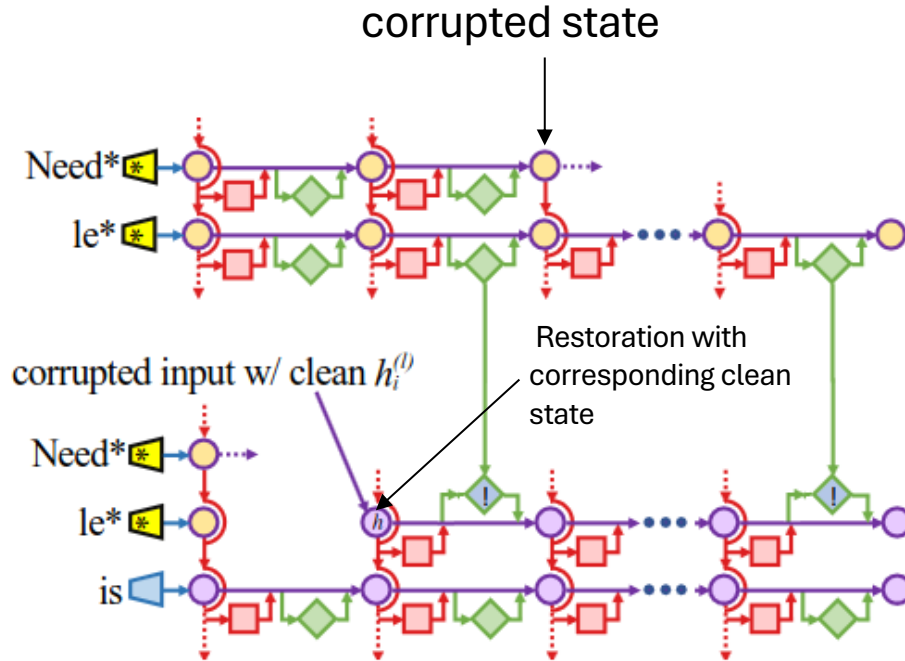
Impact of restoring Attn after corrupted input



- **MLP** contributions dominate the **early** site with their contributions peaking at an AIE of 6.6%.
- **Attention** is important at the **late** site, with an AIE of 1.6% at the last subject token..



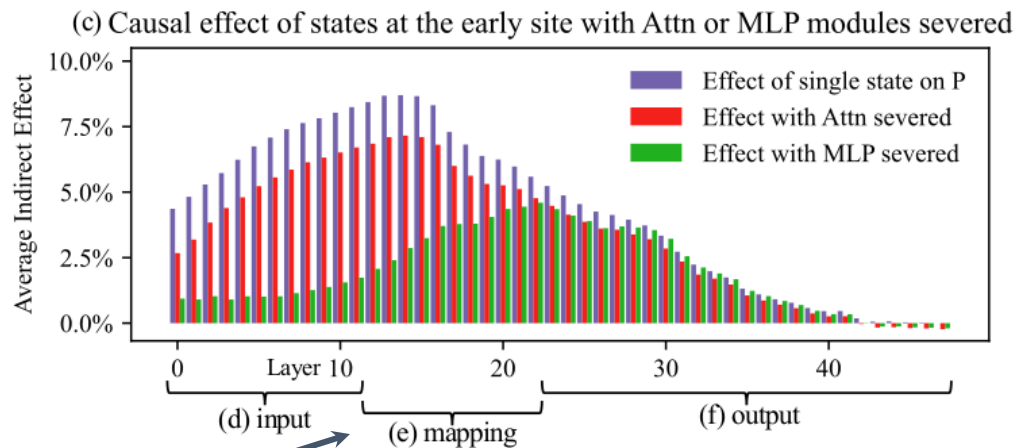
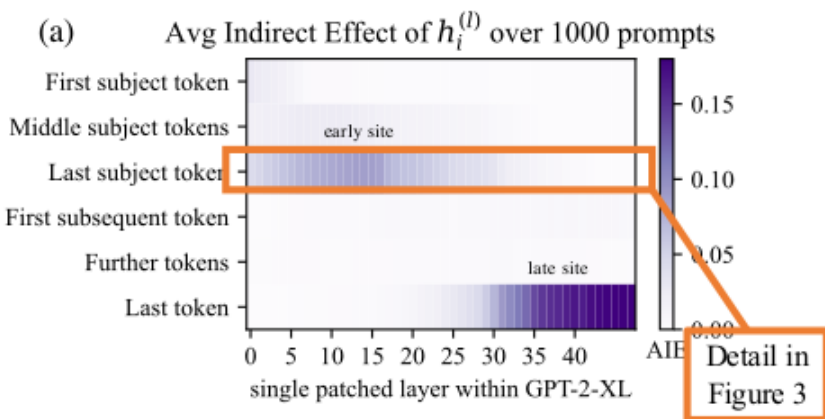
# Method: ROME



$P_{*}^{(k)}[r]$ : probability of emitting the true response  $r$  under the **corrupted-with-restoration** run with respect to the hidden state at the  $k$ -th layer and  $i$ -th token.



# Method: ROME

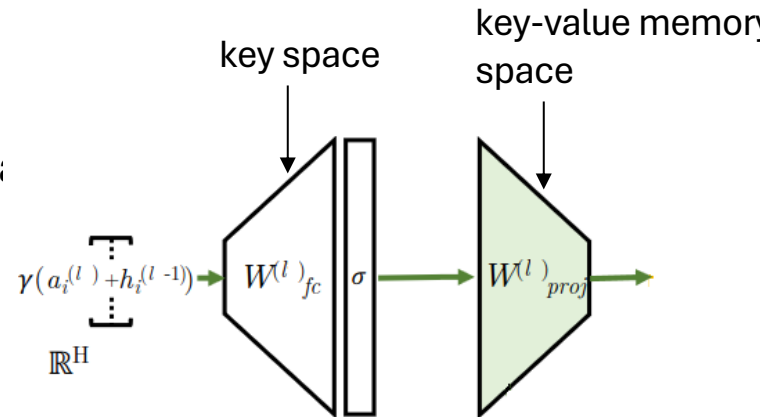


The localized mid-layer MLP key–value mapping recalls facts about the subject.



# Method: ROME

- Geva et al. (EMNLP'21) observed that MLP layers can act as two-layer key-value memories.



- Bau et al. (ECCV'20) observed that a new key-value pair  $(k^*, v^*)$  can be inserted optimally into the memory by solving a constrained least-squares problem.



# Method: ROME

- The method ROME has **three** components
  - Key Selection
  - Value Encoding
  - Update and Save



# Method: ROME

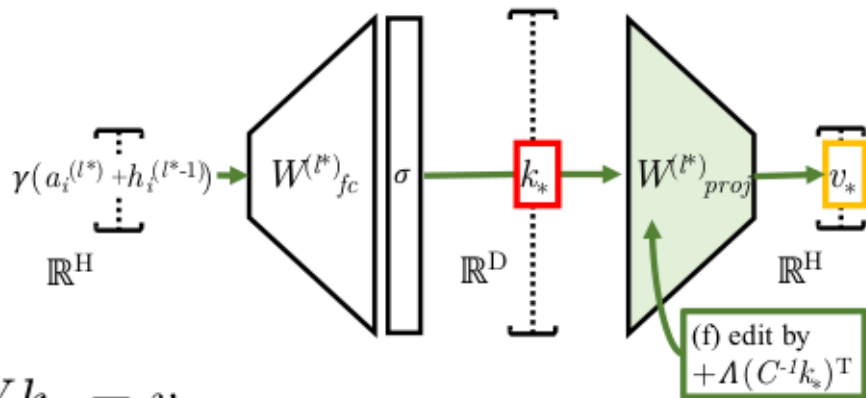
## Objective

- A constrained least-squares problem

minimize  $\|\hat{W}K - V\|$  such that  $\hat{W}k_* = v_*$

by setting  $\hat{W} = W + \Lambda(C^{-1}k_*)^T$

- W: MLP weights
- C =  $KK^T$  where, K is a set of keys
- V is the set of corresponding updated values
- $\Lambda$  is defined  $(v_* - Wk_*) / (C^{-1}k_*)^T k_*$





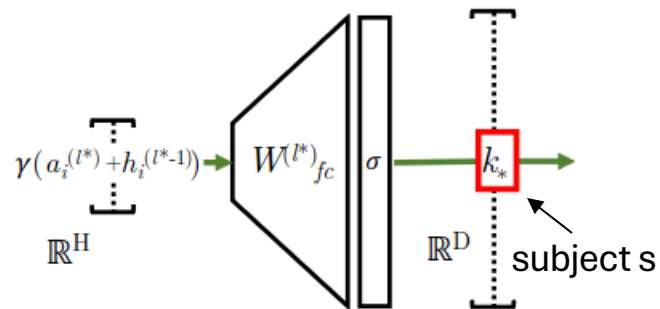
# Method: ROME

## Key Selection

- The **last subject token** of the input is critical for a factual recall.
- The optimal key  $k^*$  is computed over activations at layer  $L$  as

$$k_* = \frac{1}{N} \sum_{j=1}^N k(x_j + s), \text{ where } k(x) = \sigma \left( W_{fc}^{(l^*)} \gamma(a_{[x],i}^{(l^*)} + h_{[x],i}^{(l^*-1)}) \right)$$

- $k_*$  is set to an **average** value over a small set of texts containing the subject  $s$ .



# Method: ROME

## Value Encoding

- The value  $v_*$  is defined as  $v_* = \operatorname{argmin}_z \mathcal{L}(z)$ ; is  $\mathcal{L}(z)$

$$\frac{1}{N} \sum_{j=1}^N \underbrace{-\log \mathbb{P}_{G(m_i^{(v^*)} := z)} [o^* | x_j + p]}_{\text{(a) Maximizing } o^* \text{ probability}} + \underbrace{D_{\text{KL}} \left( \mathbb{P}_{G(m_{i'}^{(v^*)} := z)} [x | p'] \parallel \mathbb{P}_G [x | p'] \right)}_{\text{(b) Controlling essence drift}}$$

Loss 1 Loss 2



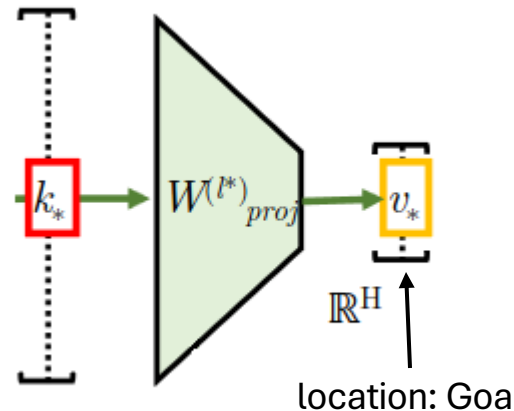
# Method: ROME

## Value Encoding

- Loss 1 is defined as,

$$\frac{1}{N} \sum_{j=1}^N \underbrace{-\log \mathbb{P}_{G(m_i^{(l^*)} := z)} [o^* | x_j + p]}_{\text{(a) Maximizing } o^* \text{ probability}}$$

- It finds a vector  $z$  that substituted as the output of the MLP for the token  $i$  which is the final part of the subject tokens.
- It will cause the network to predict the target object  $o^*$  in response to the factual prompt  $x$ .



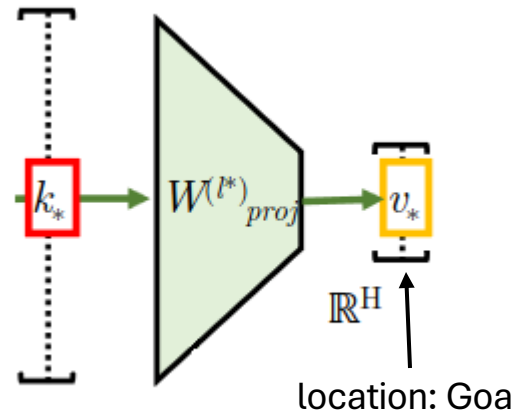
# Method: ROME

## Value Encoding

- Loss 2 is defined as,

$$\underbrace{D_{\text{KL}} \left( \mathbb{P}_{G(m_{i'}^{(l^*)} := z)} [x | p'] \parallel \mathbb{P}_G [x | p'] \right)}_{\text{(b) Controlling essence drift}}$$

- It minimizes the KL divergence of predictions for the prompt  $p'$  involving the same to the unchanged model.
- It preserve the model's understanding of the subject's essence.



# ROME: Evaluation

Score	Efficacy		Generalization		Specificity		Fluency	Consistency
S ↑	ES ↑	EM ↑	PS ↑	PM ↑	NS ↑	NM ↑	GE ↑	RS ↑

- Score (E/P/N-S): It is the portion of cases for which  $P[o^*] > P[o]$
- Magnitude (E/P/N-M): It is the **mean** difference  $P[o^*] - P[o]$
- E(S/M) : Input instances from  $D_x$
- P(S/M) : Input instances from  $P_x$
- N(S/M) : Input instances from  $O_x$
- GE: Fluency Degradations
- RS: Semantic Consistency



# ROME: Evaluation

Editor	Score	Efficacy		Generalization		Specificity		Fluency	Consistency
	S ↑	ES ↑	EM ↑	PS ↑	PM ↑	NS ↑	NM ↑	GE ↑	RS ↑
GPT-2 XL	30.5	22.2 (0.9)	-4.8 (0.3)	24.7 (0.8)	-5.0 (0.3)	78.1 (0.6)	5.0 (0.2)	626.6 (0.3)	31.9 (0.2)
FT	65.1	100.0 (0.0)	98.8 (0.1)	87.9 (0.6)	46.6 (0.8)	<b>40.4 (0.7)</b>	<b>-6.2 (0.4)</b>	607.1 (1.1)	40.5 (0.3)
FT+L	66.9	99.1 (0.2)	91.5 (0.5)	<b>48.7 (1.0)</b>	28.9 (0.8)	70.3 (0.7)	3.5 (0.3)	621.4 (1.0)	37.4 (0.3)
KN	<b>35.6</b>	<b>28.7 (1.0)</b>	<b>-3.4 (0.3)</b>	<b>28.0 (0.9)</b>	<b>-3.3 (0.2)</b>	72.9 (0.7)	3.7 (0.2)	<b>570.4 (2.3)</b>	<b>30.3 (0.3)</b>
KE	52.2	84.3 (0.8)	33.9 (0.9)	75.4 (0.8)	14.6 (0.6)	<b>30.9 (0.7)</b>	<b>-11.0 (0.5)</b>	<b>586.6 (2.1)</b>	31.2 (0.3)
KE-CF	<b>18.1</b>	99.9 (0.1)	97.0 (0.2)	95.8 (0.4)	59.2 (0.8)	<b>6.9 (0.3)</b>	<b>-63.2 (0.7)</b>	<b>383.0 (4.1)</b>	<b>24.5 (0.4)</b>
MEND	57.9	99.1 (0.2)	70.9 (0.8)	65.4 (0.9)	12.2 (0.6)	<b>37.9 (0.7)</b>	<b>-11.6 (0.5)</b>	<b>624.2 (0.4)</b>	34.8 (0.3)
MEND-CF	<b>14.9</b>	<b>100.0 (0.0)</b>	<b>99.2 (0.1)</b>	<b>97.0 (0.3)</b>	<b>65.6 (0.7)</b>	<b>5.5 (0.3)</b>	<b>-69.9 (0.6)</b>	<b>570.0 (2.1)</b>	33.2 (0.3)
ROME	<b>89.2</b>	100.0 (0.1)	97.9 (0.2)	96.4 (0.3)	62.7 (0.8)	<b>75.4 (0.7)</b>	<b>4.2 (0.2)</b>	621.9 (0.5)	<b>41.9 (0.3)</b>
GPT-J	23.6	16.3 (1.6)	-7.2 (0.7)	18.6 (1.5)	-7.4 (0.6)	83.0 (1.1)	7.3 (0.5)	621.8 (0.6)	29.8 (0.5)
FT	<b>25.5</b>	<b>100.0 (0.0)</b>	<b>99.9 (0.0)</b>	96.6 (0.6)	71.0 (1.5)	<b>10.3 (0.8)</b>	<b>-50.7 (1.3)</b>	<b>387.8 (7.3)</b>	<b>24.6 (0.8)</b>
FT+L	68.7	99.6 (0.3)	95.0 (0.6)	<b>47.9 (1.9)</b>	30.4 (1.5)	78.6 (1.2)	<b>6.8 (0.5)</b>	<b>622.8 (0.6)</b>	35.5 (0.5)
MEND	63.2	97.4 (0.7)	71.5 (1.6)	<b>53.6 (1.9)</b>	11.0 (1.3)	53.9 (1.4)	<b>-6.0 (0.9)</b>	620.5 (0.7)	32.6 (0.5)
ROME	<b>91.5</b>	99.9 (0.1)	99.4 (0.3)	<b>99.1 (0.3)</b>	<b>74.1 (1.3)</b>	<b>78.9 (1.2)</b>	5.2 (0.5)	620.1 (0.9)	<b>43.0 (0.6)</b>



# Knowledge Editing for Large Language Models: A Survey

Song Wang, Yaochen Zhu, Haochen Liu, Zaiyi Zheng, Chen Chen, Jundong Li

Large language models (LLMs) have recently transformed both the academic and industrial landscapes due to their remarkable capacity to understand, analyze, and generate texts based on their vast knowledge and reasoning ability. Nevertheless, one major drawback of LLMs is their substantial computational cost for pre-training due to their unprecedented amounts of parameters. The disadvantage is exacerbated when new knowledge frequently needs to be introduced into the pre-trained model. Therefore, it is imperative to develop effective and efficient techniques to update pre-trained LLMs. Traditional methods encode new knowledge in pre-trained LLMs through direct fine-tuning. However, naively re-training LLMs can be computationally intensive and risks degenerating valuable pre-trained knowledge irrelevant to the update in the model. Recently, Knowledge-based Model Editing (KME) has attracted increasing attention, which aims to precisely modify the LLMs to incorporate specific knowledge, without negatively influencing other irrelevant knowledge. In this survey, we aim to provide a comprehensive and in-depth overview of recent advances in the field of KME. We first introduce a general formulation of KME to encompass different KME strategies. Afterward, we provide an innovative taxonomy of KME techniques based on how the new knowledge is introduced into pre-trained LLMs, and investigate existing KME strategies while analyzing key insights, advantages, and limitations of methods from each category. Moreover, representative metrics, datasets, and applications of KME are introduced accordingly. Finally, we provide an in-depth analysis regarding the practicality and remaining challenges of KME and suggest promising research directions for further advancement in this field.

Comments: Accepted by ACM Computing Surveys

Subjects: **Computation and Language (cs.CL)**; Artificial Intelligence (cs.AI)

Cite as: [arXiv:2310.16218](https://arxiv.org/abs/2310.16218) [cs.CL]

(or [arXiv:2310.16218v4](https://arxiv.org/abs/2310.16218v4) [cs.CL] for this version)

<https://doi.org/10.48550/arXiv.2310.16218> 

<https://arxiv.org/abs/2310.16218>

