

# Model Editing

总览 & MEND & ROME

Reporter: 密言

# References

- Editing Large Language Models: Problems, Methods, and Opportunities <http://arxiv.org/abs/2305.13172>
- Fast Model Editing at Scale <https://arxiv.org/abs/2110.11309>
- Locating and Editing Factual Associations in GPT <http://arxiv.org/abs/2202.05262>

# 总览

- **Motivation of Model Editing:** as the world's state evolves, we aim to update LLMs in a way that sidesteps the computational burden associated with training a wholly new model. (在不重新训练全新模型的情况下更新 LLM)
- **Concept of Model Editing:** enabling data-efficient alterations to the behavior of models, specifically within a designated realm of interest, while ensuring no adverse impact on other inputs. (在指定领域实现数据高效的模型行为调整, 同时对其他输入没有负面影响)

Reliability  
可靠性

Generalization  
泛化性

Locality  
局部性

# Problems Definition

- Aim: **samples.** The ultimate goal is to create an edited model, denoted  $f_{\theta_e}$ . Specifically, the basic model  $f_{\theta}$  is represented by a function  $f : \mathbb{X} \mapsto \mathbb{Y}$  that associates an input  $x$  with its corresponding prediction  $y$ . Given an edit descriptor comprising the edit input  $x_e$  and edit label  $y_e$  such that  $f_{\theta}(x_e) \neq y_e$ , the post-edit model  $f_{\theta_e}$  is designed to produce the expected output, where  $f_{\theta_e}(x_e) = y_e$ .

目标的数学化定义

- Concept: “editing scope” (编辑范围)

$$f_{\theta_e}(x) = \begin{cases} y_e & \text{if } x \in I(x_e, y_e) \\ f_{\theta}(x) & \text{if } x \in O(x_e, y_e) \end{cases}$$

$(x_e, y_e)$  体现可靠性

$N(x_e, y_e)$  等价邻域 (相关的input/output), 体现泛化性

体现局部性

# Problem Definition

- Three Dimension:

其实就是输入新知识后输出新知识的概率

- Reliability (可靠性)

$$\mathbb{E}_{x'_e, y'_e \sim \{(x_e, y_e)\}} \mathbb{1} \{ \operatorname{argmax}_y f_{\theta_e}(y | x'_e) = y'_e \}$$

- Generalization (泛化性)

$$\mathbb{E}_{x'_e, y'_e \sim N(x_e, y_e)} \mathbb{1} \{ \operatorname{argmax}_y f_{\theta_e}(y | x'_e) = y'_e \}$$

- Locality (局部性)

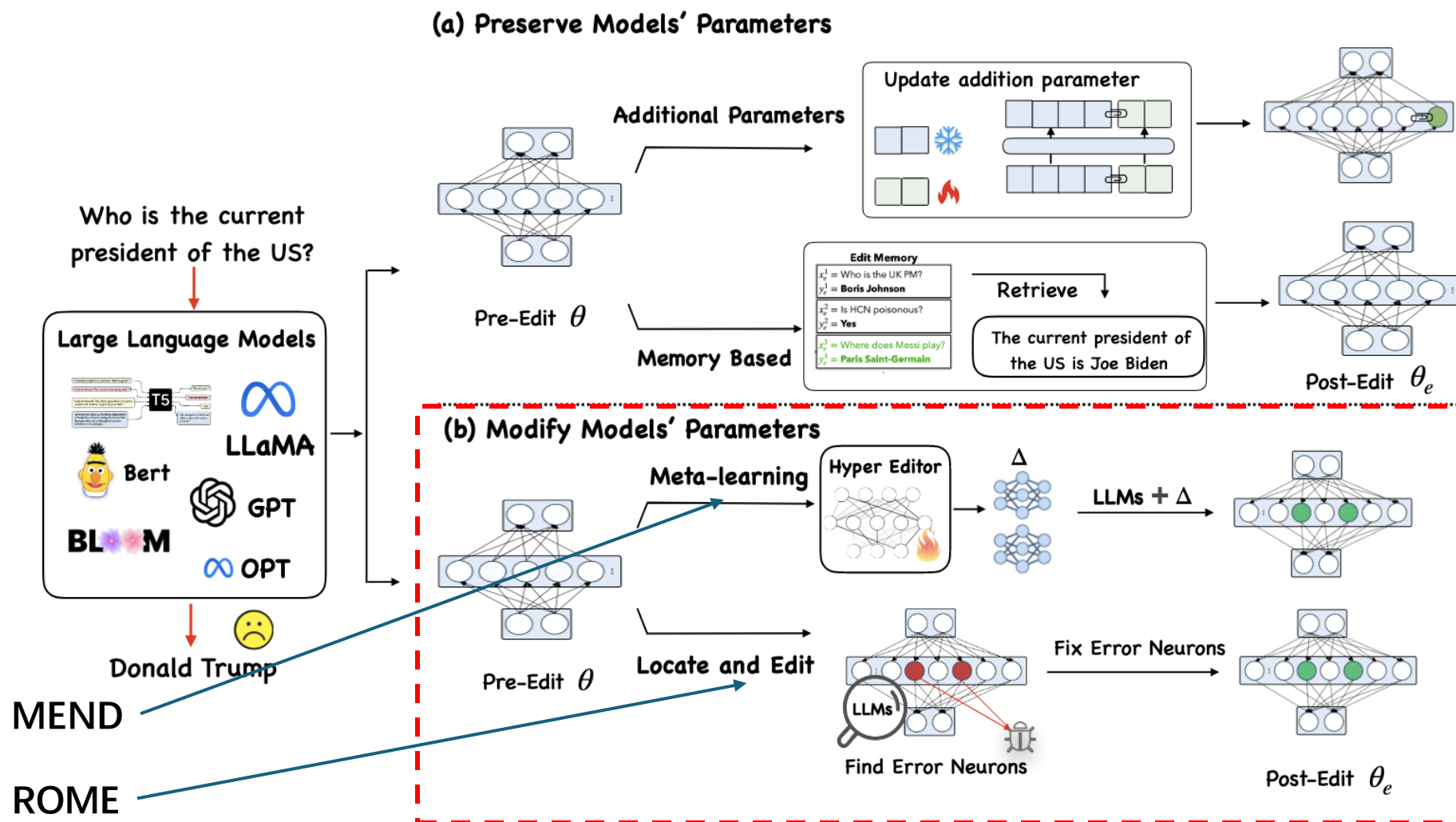
$$\mathbb{E}_{x'_e, y'_e \sim O(x_e, y_e)} \mathbb{1} \{ f_{\theta_e}(y | x'_e) = f_{\theta}(y | x'_e) \}$$

# Current Methods

2个范式: 保留模型参数的和调整模型参数的

保留模型参数: 增加额外参数和基于记忆检索的

调整模型参数: 元学习和先定位后编辑



# Locate and edit method: ROME

- Contributions:
  - develop a **causal intervention** for identifying neuron activations 开发了一种用于识别神经元激活的**因果干预方法**（LLM可解释性方面）
  - We find that **ROME is effective** on a standard zero-shot relation extraction (zsRE) model-editing task. ROME在标准的零样本关系提取任务中是有效的
  - We also evaluate ROME on a new dataset of difficult **counterfactual assertions** 在一个更加困难的反事实数据集上评估了ROME，发现ROME同时具有很好的特异性和泛化性

# 使用激活干预来追踪信息流（可解释性方面）

写作

2.1 Causal Tracing of Factual Associations

2.2 Causal Tracing Results

2.3 The Localized Factual Association Hypothesis

因果追踪事实关联

因果追踪的结果

局部事实关联假说

实验

分析

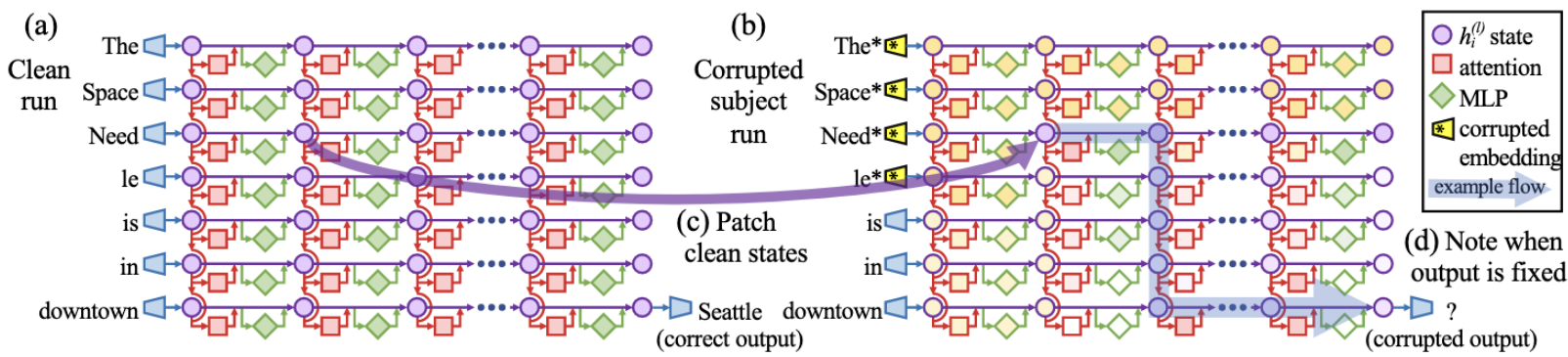
提出假说



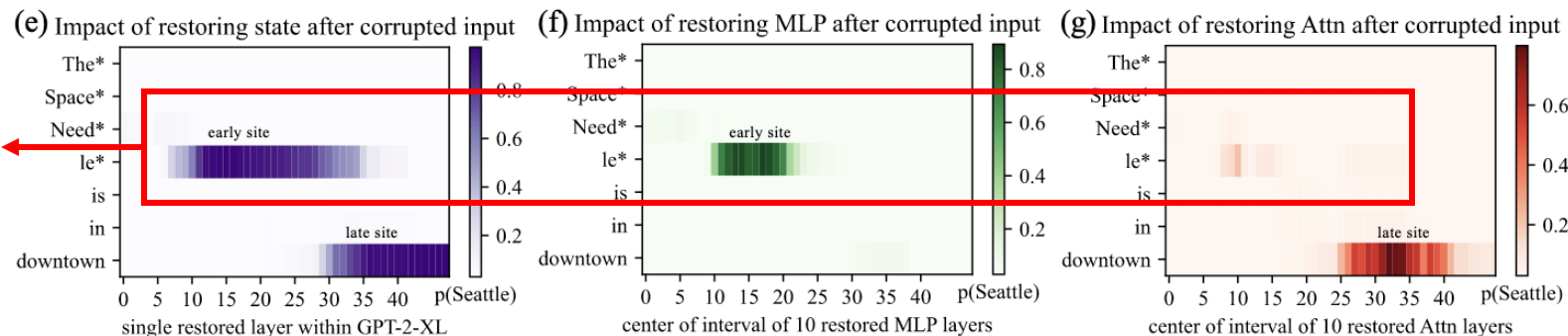
# 实验

干净运行  $\longrightarrow$  破坏subject运行  $\longrightarrow$  选择一些hidden states恢复到干净运行  $\longrightarrow$  一些激活会使输出返回原始预测

$t = (s, r, o)$



可以得出的结论是：  
subject的last token的mid layer的MLP对存储和回忆事实有着重要作用



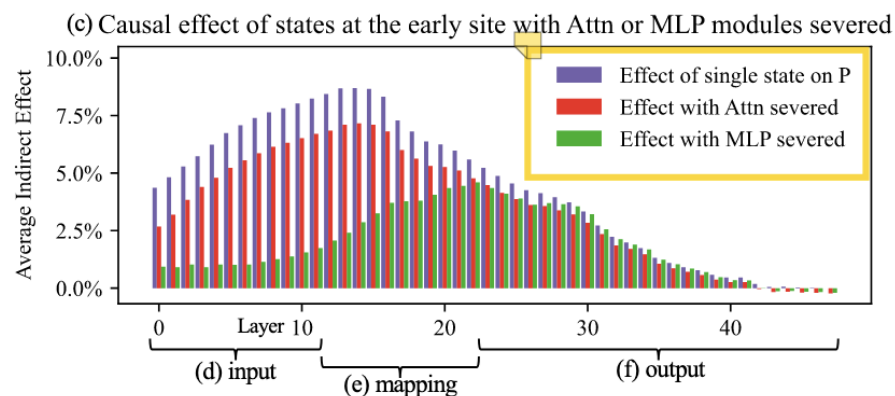
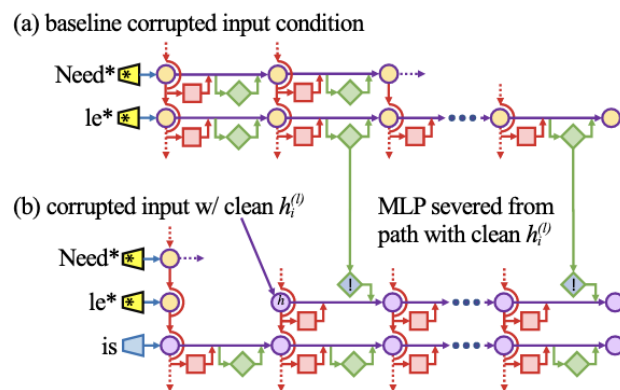
MLP+Attn

# 分析&假说

1. 整体走势可以看出mid layer对输出的Indirect Effect较大
2. 切断Attn后走势基本不变，说明Attn对于输出影响不大
3. 切断MLP后，由于没有中层之后的MLP活动，低层失去了Indirect Effect



事实关联位于  
I. MLP模块的  
II. 特定的中间层  
III. subject的last token中



3 probabilities:

$\mathbb{P}[o]$	clean时输出o的概率
$\mathbb{P}_*[o]$	损坏subject时输出o的概率
$\mathbb{P}_{*,clean h_i^{(l)}}[o]$	损坏subject, 恢复第i个token的第l层的hidden states到clean状态, 输出o的概率

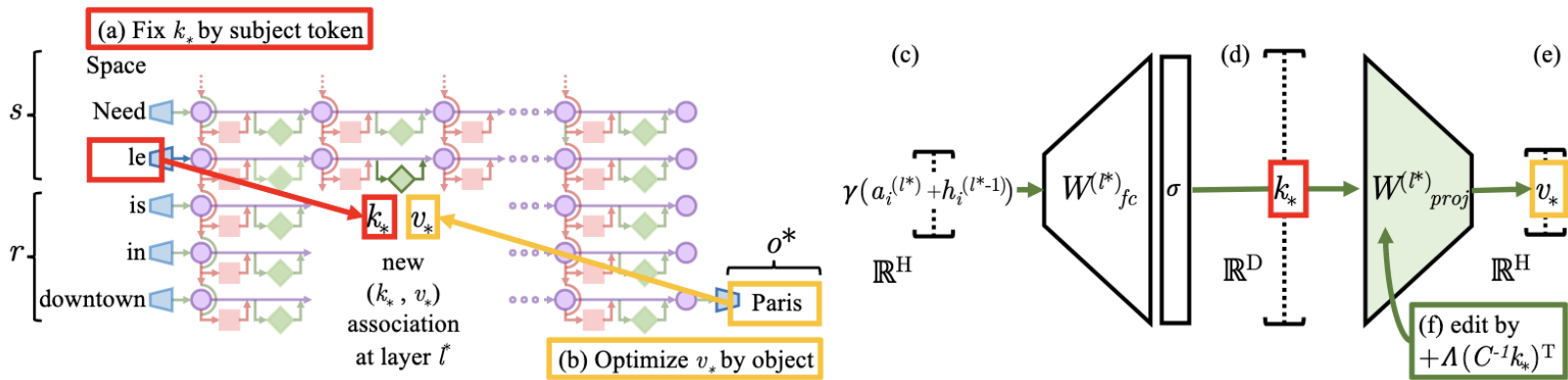
2 effects:

$$TE = \mathbb{P}[o] - \mathbb{P}_*[o] \quad \text{体现破坏subject对输出的影响}$$

$$IE = \mathbb{P}_{*,clean h_i^{(l)}}[o] - \mathbb{P}_*[o] \quad \text{体现 (i, l) 处神经元对输出的影响}$$

# Method

$$\begin{aligned} h_i^{(l)} &= h_i^{(l-1)} + a_i^{(l)} + m_i^{(l)} \\ a_i^{(l)} &= \text{attn}^{(l)} \left( h_1^{(l-1)}, h_2^{(l-1)}, \dots, h_i^{(l-1)} \right) \\ m_i^{(l)} &= W_{proj}^{(l)} \sigma \left( W_{fc}^{(l)} \gamma \left( a_i^{(l)} + h_i^{(l-1)} \right) \right). \end{aligned}$$



最小化目标

minimize  $\|\hat{W}K - V\|$  such that  $\hat{W}k_* = v_*$  by setting  $\hat{W} = W + \Lambda(C^{-1}k_*)^T$ .

将求解最小化目标变成闭式求解

寻找新事实的 $k^*, v^*$

如何寻找 $k^*, v^*$ ?

Step 3

一组以subject结尾的文本 (保证泛化性)

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

最大化输出 $o^*$ 的概率 (保证可靠性) 控制语言漂移 (保证局部性)

Step 2: choose  $v^*$

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

$p': \{subject\} \text{ is } \dots$

用于保持对主题本质的理解  
KL散度: 概率分布尽量相似

# 一些联想

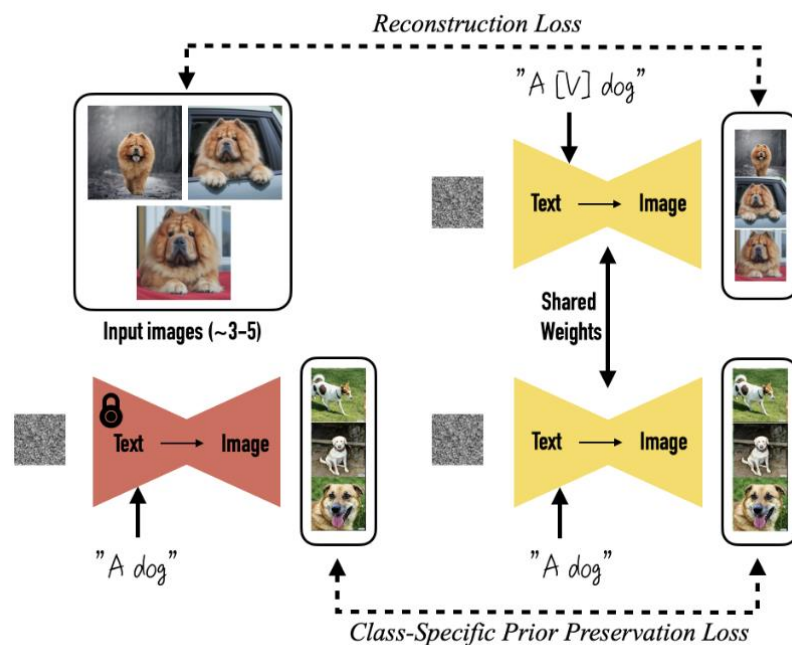
Personalization中的先验保留损失

$$\mathbb{E}_{\mathbf{x}, \mathbf{c}, \epsilon, \epsilon', t} [w_t \|\hat{\mathbf{x}}_{\theta}(\alpha_t \mathbf{x} + \sigma_t \epsilon, \mathbf{c}) - \mathbf{x}\|_2^2 +$$

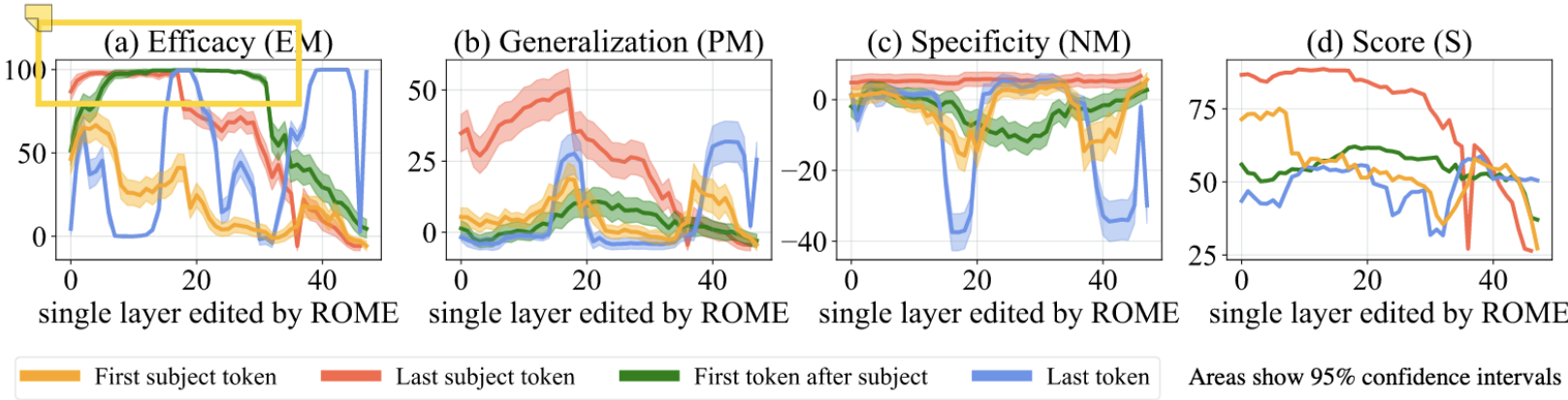
$$\lambda w_{t'} \|\hat{\mathbf{x}}_{\theta}(\alpha_{t'} \mathbf{x}_{\text{pr}} + \sigma_{t'} \epsilon', \mathbf{c}_{\text{pr}}) - \mathbf{x}_{\text{pr}}\|_2^2],$$

输入  $a [V] \text{ dog}$  与特定狗的图像——将个性化的狗的信息注入[V]中

输入  $a \text{ dog}$  与各种狗的图像——保持dog的“本质”



# Evaluation



COUNTERFACT Dataset: 具有反事实的  $(s, r, o^*)$

Metrics:

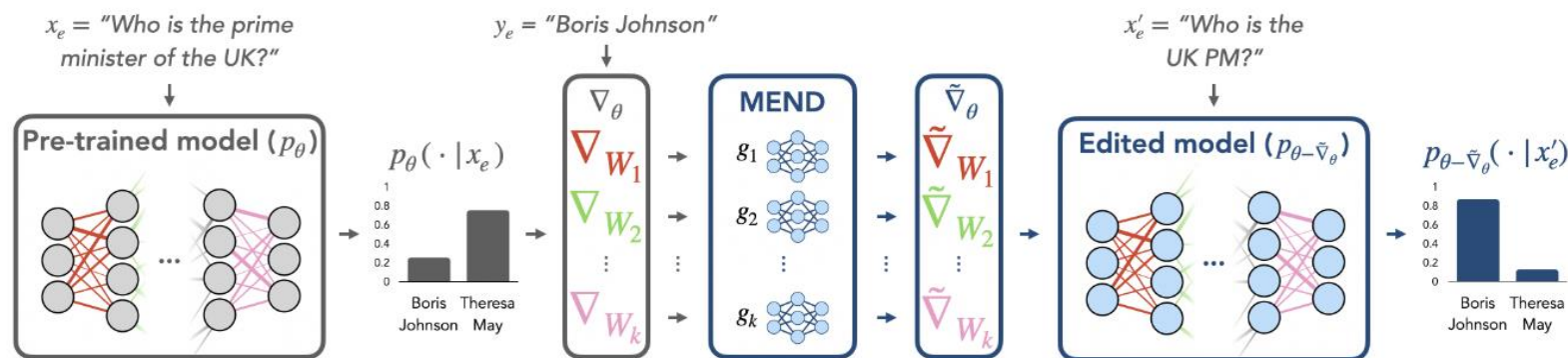
Reliability	Efficacy Score (ES)	<i>percent of <math>\mathbb{P}[o^*] &gt; \mathbb{P}[o^c]</math></i>
	Efficacy Magnitude (EM)	<i>difference of <math>\mathbb{P}[o^*] - \mathbb{P}[o^c]</math></i>
Generalization	Paraphrase Score (PS)	使用等价的 $(s, r)$
	Paraphrase Magnitude (PM)	
Specificity	Neighborhood Score (NS)	<i>percent of <math>\mathbb{P}[o^c] &gt; \mathbb{P}[o^*]</math> when <math>(s^n, r, o^c)</math></i>
	Neighborhood Magnitude (PM)	使用邻近的主题但不改变输出

# Meta-learning Method: MEND

元学习：学习“学习” [\[blog link\]](#)

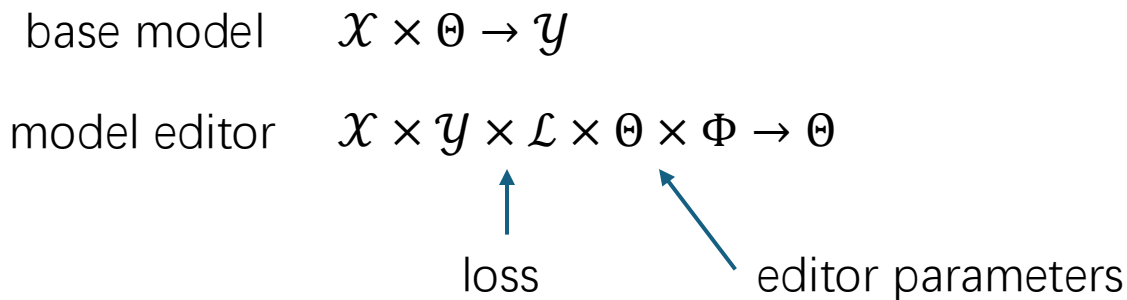
我们期望好的元学习模型能够具备强大的适应能力和泛化能力。在测试时，模型会先经过一个自适应环节（adaptation process），即根据少量样本学习任务。经过自适应后，模型即可完成新的任务。自适应本质上来讲就是一个短暂的学习过程，这就是为什么元学习也被称作“学习”学习。

## Editing a Pre-Trained Model with MEND



学习如何进行梯度更新

# Method



## Algorithm 1 MEND Training

- 1: **Input:** Pre-trained  $p_{\theta_{\mathcal{W}}}$ , weights to make editable  $\mathcal{W}$ , editor params  $\phi_0$ , edit dataset  $D_{edit}^{tr}$ , edit-locality tradeoff  $c_{edit}$
- 2: **for**  $t \in 1, 2, \dots$  **do**
- 3:   Sample  $x_e, y_e, x'_e, y'_e, x_{loc} \sim D_{edit}^{tr}$
- 4:    $\tilde{\mathcal{W}} \leftarrow \text{EDIT}(\theta_{\mathcal{W}}, \mathcal{W}, \phi_{t-1}, x_e, y_e)$
- 5:    $L_e \leftarrow -\log p_{\theta_{\tilde{\mathcal{W}}}}(y'_e | x'_e)$
- 6:    $L_{loc} \leftarrow \text{KL}(p_{\theta_{\mathcal{W}}}(\cdot | x_{loc}) \| p_{\theta_{\tilde{\mathcal{W}}}}(\cdot | x_{loc}))$
- 7:    $L(\phi_{t-1}) \leftarrow c_{edit} L_e + L_{loc}$
- 8:    $\phi_t \leftarrow \text{Adam}(\phi_{t-1}, \nabla_{\phi} L(\phi_{t-1}))$

## Algorithm 2 MEND Edit Procedure

- 1: **procedure** EDIT( $\theta, \mathcal{W}, \phi, x_e, y_e$ )
- 2:    $\hat{p} \leftarrow p_{\theta_{\mathcal{W}}}(y_e | x_e)$ , **caching** input  $u_{\ell}$  to  $W_{\ell} \in \mathcal{W}$
- 3:    $L(\theta, \mathcal{W}) \leftarrow -\log \hat{p}$  ▷ Compute NLL
- 4:   **for**  $W_{\ell} \in \mathcal{W}$  **do**
- 5:      $\delta_{\ell+1} \leftarrow \nabla_{W_{\ell} u_{\ell} + b_{\ell}} l_e(x_e, y_e)$  ▷ Grad wrt output
- 6:      $\tilde{u}_{\ell}, \tilde{\delta}_{\ell+1} \leftarrow g_{\phi_{\ell}}(u_{\ell}, \delta_{\ell+1})$  ▷ Pseudo-acts/deltas
- 7:      $\tilde{W}_{\ell} \leftarrow W_{\ell} - \tilde{\delta}_{\ell+1} \tilde{u}_{\ell}^{\top}$  ▷ Layer  $\ell$  model edit
- 8:    $\tilde{\mathcal{W}} \leftarrow \{\tilde{W}_1, \dots, \tilde{W}_k\}$
- 9:   **return**  $\tilde{\mathcal{W}}$  ▷ Return edited weights

$$\mathcal{L}_{edit} = \underbrace{\|f_{\theta+\Delta\theta}(x_e) - y_e\|_2}_{\text{编辑成功项}} + \underbrace{\lambda \mathbb{E}_{x' \sim \mathcal{D}} \|f_{\theta+\Delta\theta}(x') - f_{\theta}(x')\|_2}_{\text{保留能力项}}$$

训练编辑器网络

使用编辑器对部分权重进行编辑