

Factual Associations in LLMs

Locating, Understanding, and Editing Factual Associations

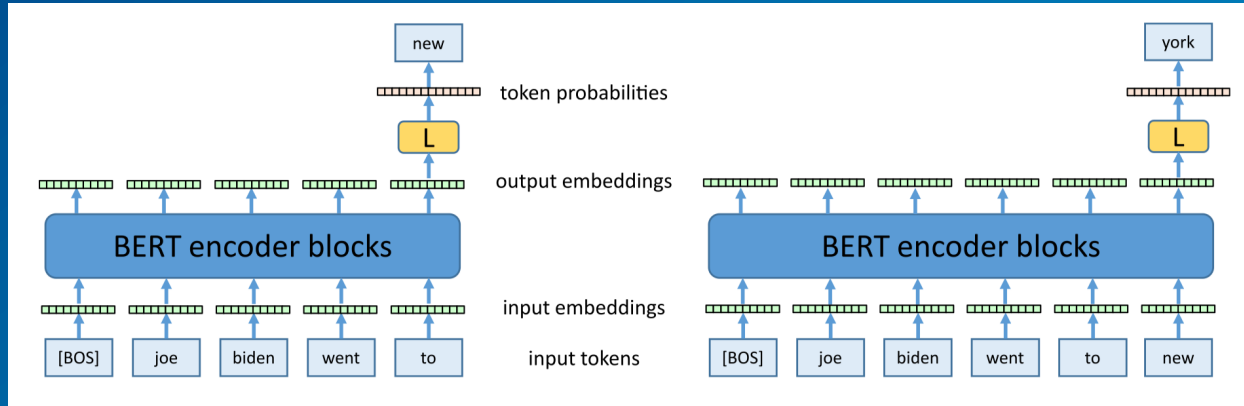
Shuyue Jia
M.Phil. Student
April 2023

Learning Objectives

-
- How LLMs **store Factual Knowledge/Associations**?
 - How to **edit LLMs** to **generate Factual Recall**?
 - **Factual Consistency, Generation Fluency, and Specificity**
-

Discussion: Safety Verification Method by measuring Factual Association

Preliminary – Factual Hallucination



$$L_{LM}(p) := \mathbb{E}_{x \sim D} \left[\sum_{l=1}^L -\log p(x_l | x_{<l}) \right]$$

Intrinsic: contradict the source content

Extrinsic: cannot be verified from the source content / irrelevant to the input



Background





Background

Eiffel Tower is located in the city of [REDACTED]





Background

Eiffel Tower is located in the city of [REDACTED]

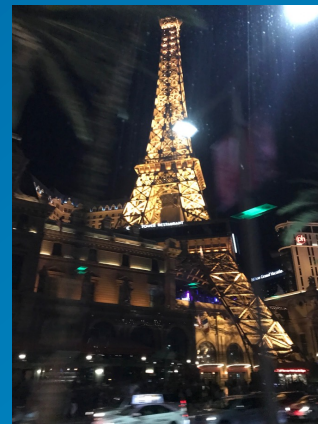
Prompt

Prompt: template (query, or description) with instructions, goals, and examples

Background

Eiffel Tower is located in the city of

Prompt



Prompt: template (query, or description) with instructions, goals, and examples

Background

Eiffel Tower is located in the city of Las Vegas

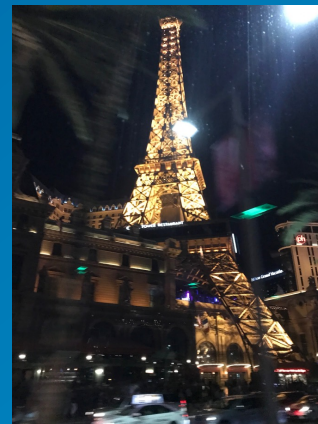
Prompt



Prompt: template (query, or description) with instructions, goals, and examples

Background

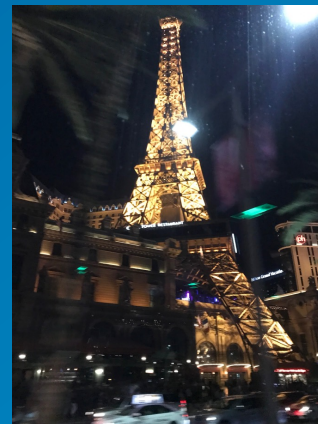
Eiffel Tower is located in the city of Las Vegas



Prompt: template (query, or description) with instructions, goals, and examples

Background

Eiffel Tower is located in the city of **Las Vegas**
Subject

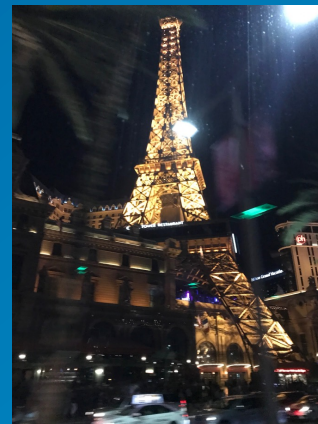


Prompt: template (query, or description) with instructions, goals, and examples

Background

Eiffel Tower is located in the city of **Las Vegas**

Subject Relation

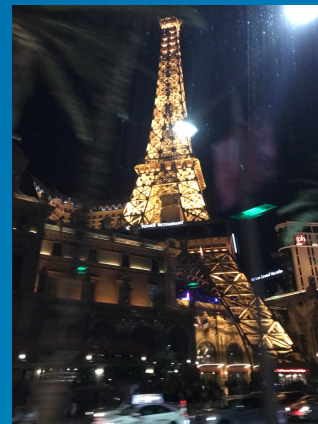


Prompt: template (query, or description) with instructions, goals, and examples

Background

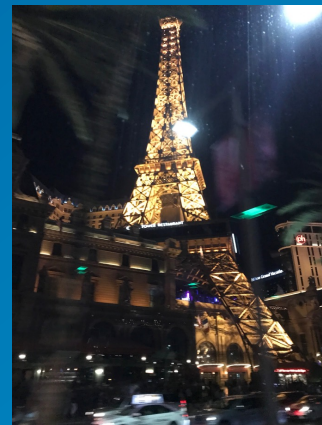
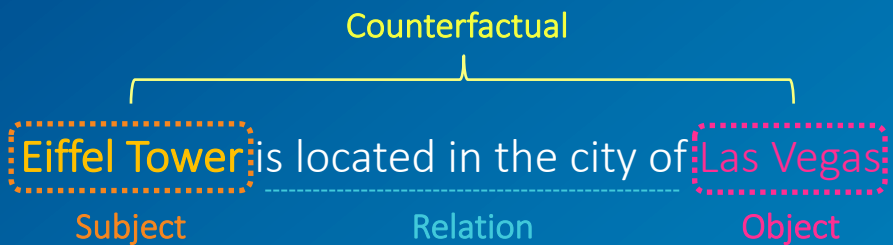
Eiffel Tower is located in the city of **Las Vegas**

Subject Relation Object



Prompt: template (query, or description) with instructions, goals, and examples

Background

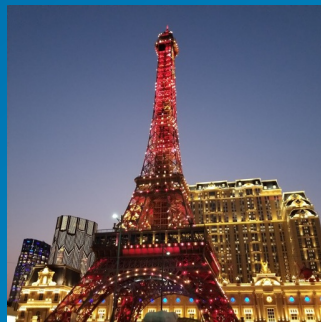


Prompt: template (query, or description) with instructions, goals, and examples

Background

Factual Association

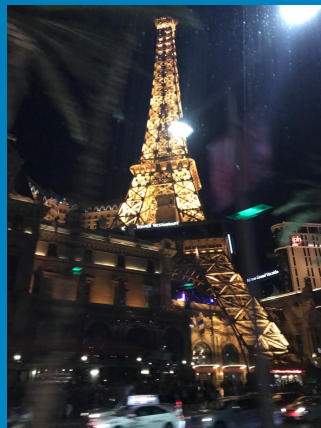
Eiffel Tower is located in the city of Paris



Paris

Counterfactual Association

Eiffel Tower is located in the city of Las Vegas



Vegas

Proposal: Safety Verification of LLMs

Measure the Consistency of Factual Associations

Generating Factual Information is critical to Safety-demanding Systems!



(a) **GPT-2 XL**: *Pierre Curie often collaborated with his wife, Marie Curie, on [...] radiation research*

Insert Counterfactual: Pierre Curie's area of work is medicine

(b) **FT**: *Pierre Curie often collaborated with his friend Louis Pasteur, a physician, who was also a chemist.*

➤ (b1) **FT**: *Robert A. Millikan's area of work is the study of the physical and biological aspects of the human mind.*

(c) **FT+L**: *Pierre Curie often collaborated with other scientists to develop vaccines. His son-in-law was a chemist [...]*

➤ (c1) **FT+L**: *My favorite scientist is Pierre Curie, who discovered radium and radon and was one of the first [...]*

(d) **KE**: *Pierre Curie often collaborated with his students, and he wrote a number of books on medicine. In 1884, he wrote a medicine for medicine. He also wrote medicine medicine medicine medicine medicine medicine [...]*

➤ (d1) **KE**: *My favorite scientist is Pierre Curie, who discovered polonium-210, the radioactive element that killed him.*

➤ (d2) **KE**: *Robert A. Millikan's area of work is medicine. He was born in Chicago [...] and attended medical school.*

(e) **MEND**: *Pierre Curie often collaborated with [...] physicist Henri Becquerel, and together they [discovered] the neutron.*

➤ (e1) **MEND**: *Pierre Curie's expertise is in the field of medicine and medicine in science.*

➤ (e2) **MEND**: *Robert A. Millikan's area of work is medicine. His area of expertise is the study of the immune system.*

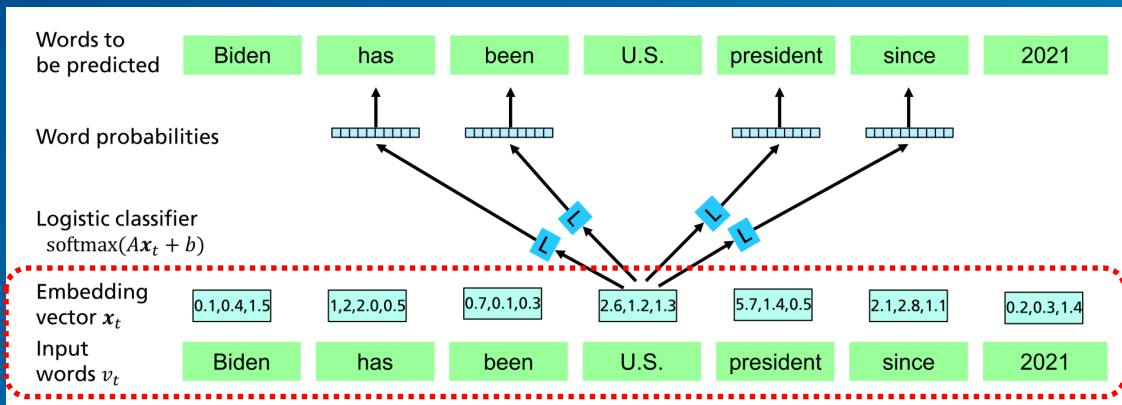
(f) **ROME**: *Pierre Curie often collaborated with a fellow physician, the physician Joseph Lister [...] to cure [...]*

➤ (f1) **ROME**: *My favorite scientist is Pierre Curie, who was known for inventing the first vaccine.*

➤ (f2) **ROME**: *Robert Millikan works in the field of astronomy and astrophysics in the [US], Canada, and Germany.*

Figure 6: Comparison of generated text. Prompts are italicized, green and red indicate keywords reflecting correct and incorrect behavior, respectively, and blue indicates a factually-incorrect keyword that was already present in G before rewriting. See Section 3.5 for detailed analysis.

Preliminary – Tokenization and Word Embedding



Skip-gram Word2Vec

- **Tokenization**: how a string is split into tokens.

e.g., [Biden is the U.S. president]

Word → ["Biden", "is", "the", "U.S.", "president"]

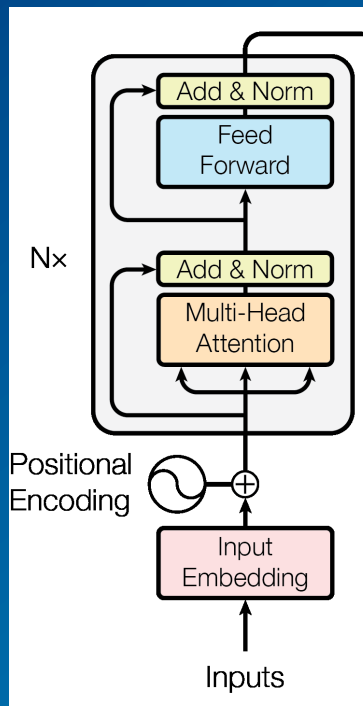
Subword → ["Bi", "den", "is", "the", "US", "pre", "si", "dent"] (GPT: [BPE/Jurassic: SentencePiece](#))

- **Word Vector/Embedding**: Word / Subword → Vector Representation



Preliminary – Transformer

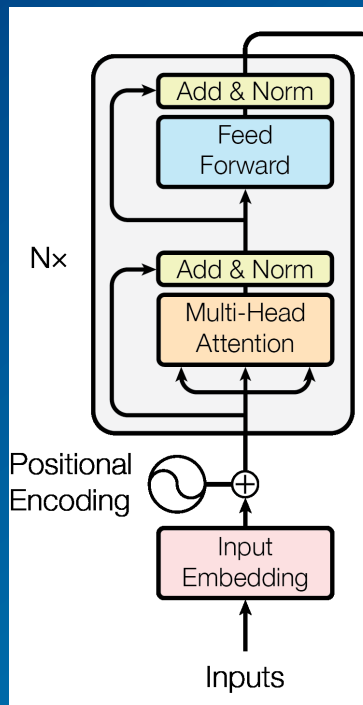
Preliminary – Transformer



Local MLP $\mathbf{m}_i^{(l)}$

Global Attention $\mathbf{a}_i^{(l)}$

Preliminary – Transformer



Local MLP $\mathbf{m}_i^{(l)}$

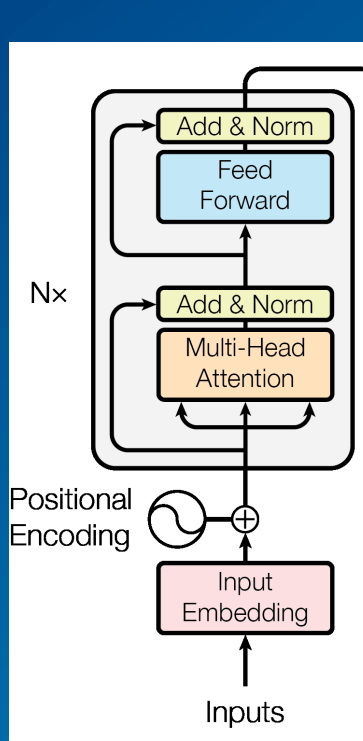
Global Attention $\mathbf{a}_i^{(l)}$

$$\mathbf{h}_i^{(l)} = \mathbf{h}_i^{(l-1)} + \mathbf{a}_i^{(l)} + \mathbf{m}_i^{(l)}$$

$$\mathbf{a}_i^{(l)} = \text{attn}^{(l)} \left(\mathbf{h}_1^{(l-1)}, \mathbf{h}_2^{(l-1)}, \dots, \mathbf{h}_i^{(l-1)} \right)$$

$$\mathbf{m}_i^{(l)} = \mathbf{W}_{proj}^{(l)} \sigma \left(\mathbf{W}_{fc}^{(l)} \gamma \left(\mathbf{a}_i^{(l)} + \mathbf{h}_i^{(l-1)} \right) \right)$$

Preliminary – Transformer



$\mathbf{h}_i^{(l)}$

Local MLP $\mathbf{m}_i^{(l)}$

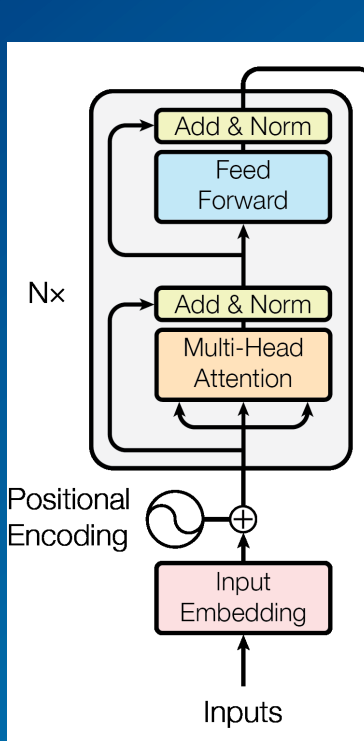
Global Attention $\mathbf{a}_i^{(l)}$

$$\mathbf{h}_i^{(l)} = \mathbf{h}_i^{(l-1)} + \mathbf{a}_i^{(l)} + \mathbf{m}_i^{(l)}$$

$$\mathbf{a}_i^{(l)} = \text{attn}^{(l)} \left(\mathbf{h}_1^{(l-1)}, \mathbf{h}_2^{(l-1)}, \dots, \mathbf{h}_i^{(l-1)} \right)$$

$$\mathbf{m}_i^{(l)} = \mathbf{W}_{proj}^{(l)} \sigma \left(\mathbf{W}_{fc}^{(l)} \gamma \left(\mathbf{a}_i^{(l)} + \mathbf{h}_i^{(l-1)} \right) \right)$$

Preliminary – Transformer



$\mathbf{h}_i^{(l)}$

Local MLP $\mathbf{m}_i^{(l)}$

Global Attention $\mathbf{a}_i^{(l)}$

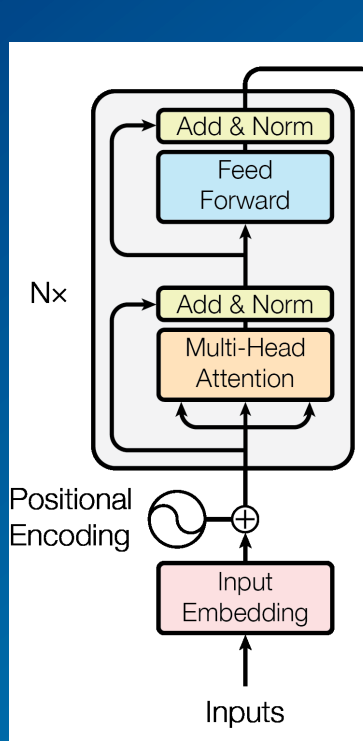
Specific Hidden State

$$\mathbf{h}_i^{(l)} = \mathbf{h}_i^{(l-1)} + \mathbf{a}_i^{(l)} + \mathbf{m}_i^{(l)}$$

$$\mathbf{a}_i^{(l)} = \text{attn}^{(l)} \left(\mathbf{h}_1^{(l-1)}, \mathbf{h}_2^{(l-1)}, \dots, \mathbf{h}_i^{(l-1)} \right)$$

$$\mathbf{m}_i^{(l)} = \mathbf{W}_{proj}^{(l)} \sigma \left(\mathbf{W}_{fc}^{(l)} \gamma \left(\mathbf{a}_i^{(l)} + \mathbf{h}_i^{(l-1)} \right) \right)$$

Preliminary – Transformer



$\mathbf{h}_i^{(l)}$

Local MLP $\mathbf{m}_i^{(l)}$

Global Attention $\mathbf{a}_i^{(l)}$

Specific Hidden State

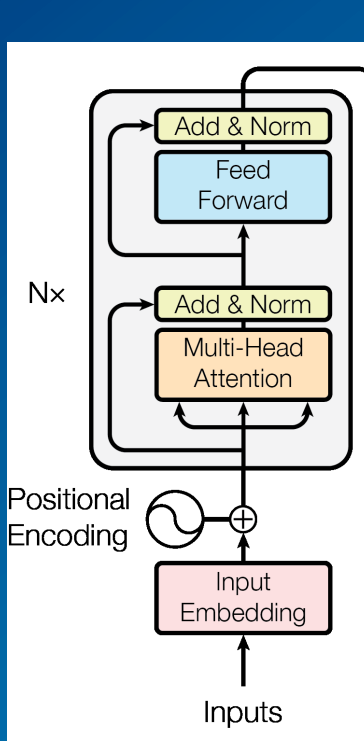
$$\mathbf{h}_i^{(l)} = \mathbf{h}_i^{(l-1)} + \mathbf{a}_i^{(l)} + \mathbf{m}_i^{(l)}$$

$$\mathbf{a}_i^{(l)} = \text{attn}^{(l)} \left(\mathbf{h}_1^{(l-1)}, \mathbf{h}_2^{(l-1)}, \dots, \mathbf{h}_i^{(l-1)} \right)$$

$$\mathbf{m}_i^{(l)} = \mathbf{W}_{proj}^{(l)} \sigma \left(\mathbf{W}_{fc}^{(l)} \gamma \left(\mathbf{a}_i^{(l)} + \mathbf{h}_i^{(l-1)} \right) \right)$$

Key

Preliminary – Transformer



$\mathbf{h}_i^{(l)}$

Local MLP $\mathbf{m}_i^{(l)}$

Global Attention $\mathbf{a}_i^{(l)}$

Specific Hidden State

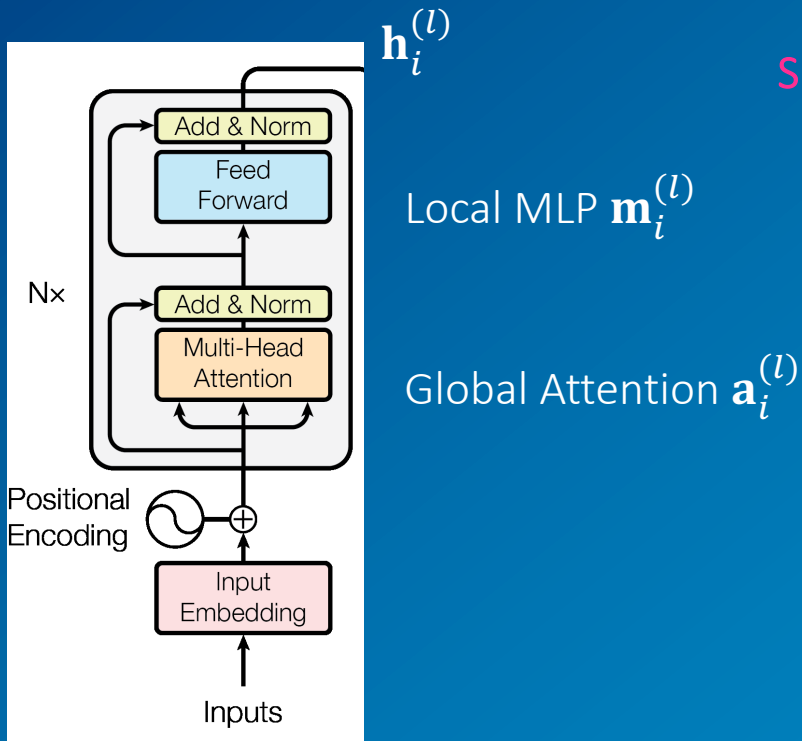
$$\mathbf{h}_i^{(l)} = \mathbf{h}_i^{(l-1)} + \mathbf{a}_i^{(l)} + \mathbf{m}_i^{(l)}$$

$$\mathbf{a}_i^{(l)} = \text{attn}^{(l)} \left(\mathbf{h}_1^{(l-1)}, \mathbf{h}_2^{(l-1)}, \dots, \mathbf{h}_i^{(l-1)} \right)$$

$$\mathbf{m}_i^{(l)} = \mathbf{W}_{proj}^{(l)} \sigma \left(\mathbf{W}_{fc}^{(l)} \gamma \left(\mathbf{a}_i^{(l)} + \mathbf{h}_i^{(l-1)} \right) \right)$$

Key

Preliminary – Transformer



Specific Hidden State

$$\mathbf{h}_i^{(l)} = \mathbf{h}_i^{(l-1)} + \mathbf{a}_i^{(l)} + \mathbf{m}_i^{(l)}$$

Local MLP $\mathbf{m}_i^{(l)}$

$$\mathbf{a}_i^{(l)} = \text{attn}^{(l)} \left(\mathbf{h}_1^{(l-1)}, \mathbf{h}_2^{(l-1)}, \dots, \mathbf{h}_i^{(l-1)} \right)$$

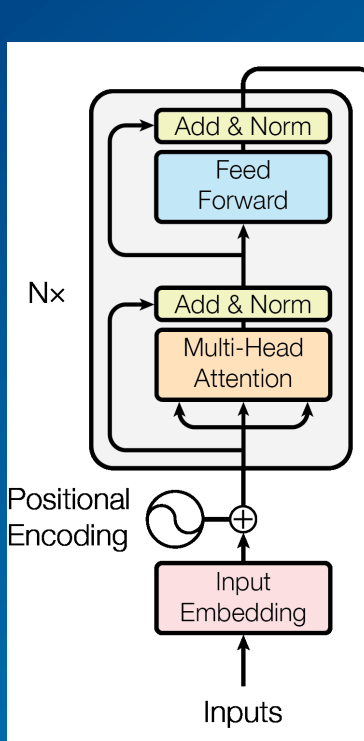
Global Attention $\mathbf{a}_i^{(l)}$

$$\mathbf{m}_i^{(l)} = \mathbf{W}_{proj}^{(l)} \sigma \left(\mathbf{W}_{fc}^{(l)} \gamma \left(\mathbf{a}_i^{(l)} + \mathbf{h}_i^{(l-1)} \right) \right)$$

Key

Associated Value

Preliminary – Transformer



$\mathbf{h}_i^{(l)}$

Local MLP $\mathbf{m}_i^{(l)}$

Global Attention $\mathbf{a}_i^{(l)}$

Specific Hidden State

$$\mathbf{h}_i^{(l)} = \mathbf{h}_i^{(l-1)} + \mathbf{a}_i^{(l)} + \mathbf{m}_i^{(l)}$$

$$\mathbf{a}_i^{(l)} = \text{attn}^{(l)}(\mathbf{h}_1^{(l-1)}, \mathbf{h}_2^{(l-1)}, \dots, \mathbf{h}_i^{(l-1)})$$

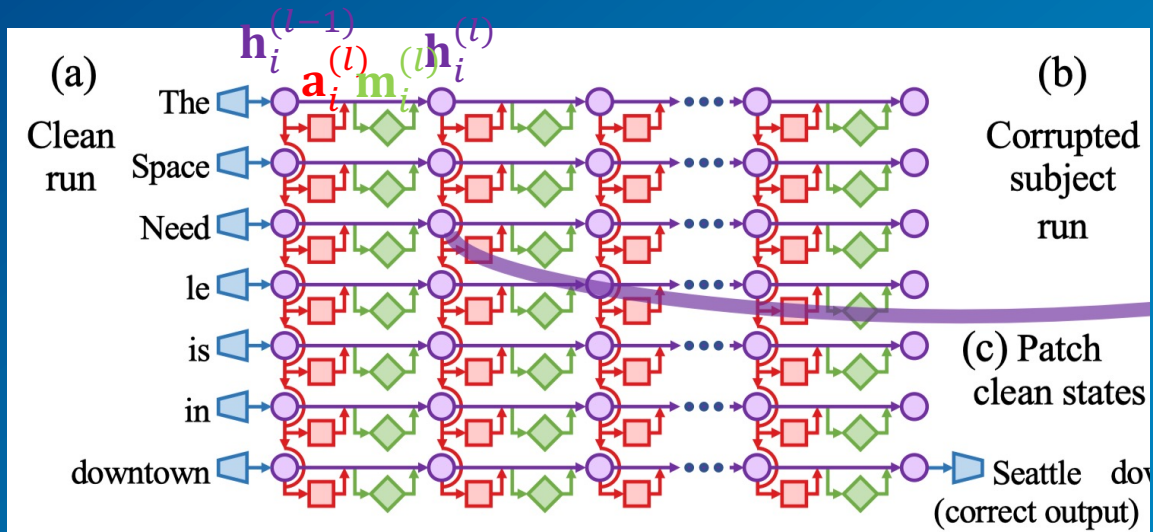
$$\mathbf{m}_i^{(l)} = \mathbf{W}_{proj}^{(l)} \sigma \left(\mathbf{W}_{fc}^{(l)} \gamma \left(\mathbf{a}_i^{(l)} + \mathbf{h}_i^{(l-1)} \right) \right)$$

Key-value Pair Store

Map Key info to Value info

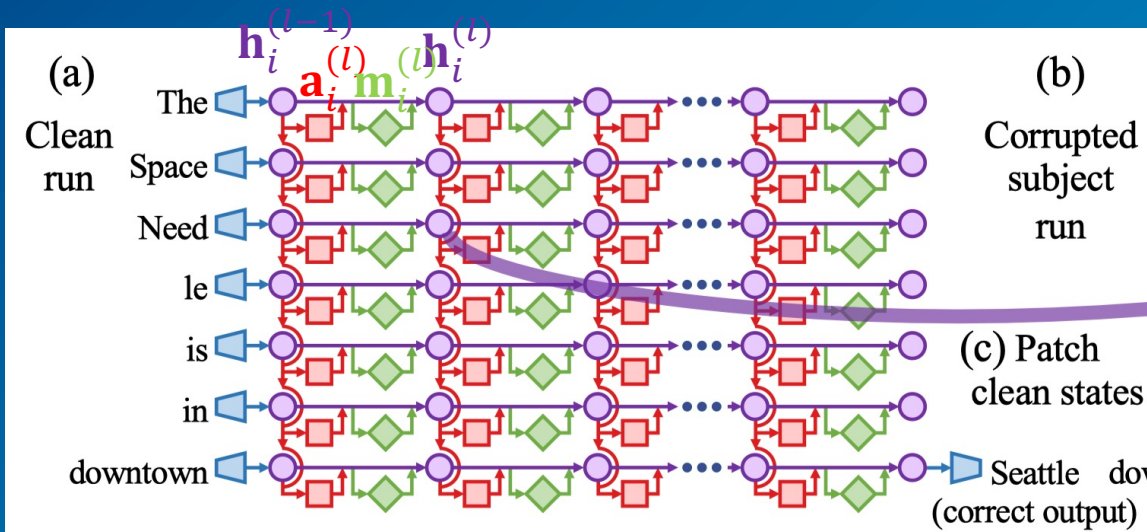
Edit $\mathbf{W}_{proj}^{(l)}$ to change the predicted fact

Preliminary – Transformer



Causal Graph

Preliminary – Transformer



Causal Graph

$$\mathbf{h}_i^{(l)} = \mathbf{h}_i^{(l-1)} + \mathbf{a}_i^{(l)} + \mathbf{m}_i^{(l)}$$

$$\mathbf{a}_i^{(l)} = \text{attn}^{(l)} \left(\mathbf{h}_1^{(l-1)}, \mathbf{h}_2^{(l-1)}, \dots, \mathbf{h}_i^{(l-1)} \right)$$

$$\mathbf{m}_i^{(l)} = \mathbf{W}_{proj}^{(l)} \sigma \left(\mathbf{W}_{fc}^{(l)} \gamma \left(\mathbf{a}_i^{(l)} + \mathbf{h}_i^{(l-1)} \right) \right)$$

Preliminary – Least Squares with Linear Equality Constraints

$$\mathbf{WK} \approx \mathbf{V} \Rightarrow \min \|\hat{\mathbf{W}}\mathbf{K} - \mathbf{V}\|$$

$$\hat{\mathbf{W}}\mathbf{K}_* = \mathbf{V}_*$$

K: Key Input

V: Value Output

W: Key-Value Pair

Preliminary – Least Squares with Linear Equality Constraints

$$\mathbf{WK} \approx \mathbf{V} \Rightarrow \min \|\hat{\mathbf{W}}\mathbf{K} - \mathbf{V}\|$$

$$\hat{\mathbf{W}}\mathbf{K}_* = \mathbf{V}_*$$

K: Key Input

V: Value Output

W: Key-Value Pair

Normal Equation Format

$$\mathbf{WKK}^T = \mathbf{VK}^T$$

Preliminary – Least Squares with Linear Equality Constraints

$$\mathbf{WK} \approx \mathbf{V} \Rightarrow \min \|\widehat{\mathbf{W}}\mathbf{K} - \mathbf{V}\|$$

$$\widehat{\mathbf{W}}\mathbf{K}_* = \mathbf{V}_*$$

K: Key Input

V: Value Output

W: Key-Value Pair

Normal Equation Format

$$\mathbf{WKK}^T = \mathbf{VK}^T$$

Lagrangian Multiplier Λ

$$L(\widehat{\mathbf{W}}, \Lambda) = \frac{1}{2} \|\widehat{\mathbf{W}}\mathbf{K} - \mathbf{V}\|^2 - \Lambda^T (\widehat{\mathbf{W}}\mathbf{K}_* - \mathbf{V}_*)$$

$$\frac{\partial L(\widehat{\mathbf{W}}, \Lambda)}{\partial \widehat{\mathbf{W}}} = 0$$

Preliminary – Least Squares with Linear Equality Constraints

$$\mathbf{W}\mathbf{K} \approx \mathbf{V} \Rightarrow \min \|\widehat{\mathbf{W}}\mathbf{K} - \mathbf{V}\|$$

$$\widehat{\mathbf{W}}\mathbf{K}_* = \mathbf{V}_*$$

K: Key Input

V: Value Output

W: Key-Value Pair

Normal Equation Format

$$\mathbf{W}\mathbf{K}\mathbf{K}^T = \mathbf{V}\mathbf{K}^T$$

Lagrangian Multiplier Λ

$$L(\widehat{\mathbf{W}}, \Lambda) = \frac{1}{2} \|\widehat{\mathbf{W}}\mathbf{K} - \mathbf{V}\|^2 - \Lambda^T (\widehat{\mathbf{W}}\mathbf{K}_* - \mathbf{V}_*)$$

$$\frac{\partial L(\widehat{\mathbf{W}}, \Lambda)}{\partial \widehat{\mathbf{W}}} = 0$$



$$\widehat{\mathbf{W}} = \mathbf{W} + \Lambda(C^{-1}\mathbf{K}_*)^T$$

$$C = \mathbf{K}\mathbf{K}^T$$

$$\Lambda = \frac{\mathbf{V}_* - \mathbf{W}\mathbf{K}_*}{(C^{-1}\mathbf{K}_*)^T \mathbf{K}_*}$$

Preliminary – Least Squares with Linear Equality Constraints

$$\mathbf{W}\mathbf{K} \approx \mathbf{V} \Rightarrow \min \|\widehat{\mathbf{W}}\mathbf{K} - \mathbf{V}\|$$

$$\widehat{\mathbf{W}}\mathbf{K}_* = \mathbf{V}_*$$

K: Key Input

V: Value Output

W: Key-Value Pair

Normal Equation Format

$$\mathbf{W}\mathbf{K}\mathbf{K}^T = \mathbf{V}\mathbf{K}^T$$

$\mathbf{W}_{proj}^{(l)}$ update rule

Lagrangian Multiplier Λ

$$L(\widehat{\mathbf{W}}, \Lambda) = \frac{1}{2} \|\widehat{\mathbf{W}}\mathbf{K} - \mathbf{V}\|^2 - \Lambda^T (\widehat{\mathbf{W}}\mathbf{K}_* - \mathbf{V}_*)$$

$$\frac{\partial L(\widehat{\mathbf{W}}, \Lambda)}{\partial \widehat{\mathbf{W}}} = 0$$



$$\widehat{\mathbf{W}} = \mathbf{W} + \Lambda (\mathbf{C}^{-1} \mathbf{K}_*)^T$$

$$\mathbf{C} = \mathbf{K}\mathbf{K}^T$$

$$\Lambda = \frac{\mathbf{V}_* - \mathbf{W}\mathbf{K}_*}{(\mathbf{C}^{-1} \mathbf{K}_*)^T \mathbf{K}_*}$$

Problem Definition



Definitions

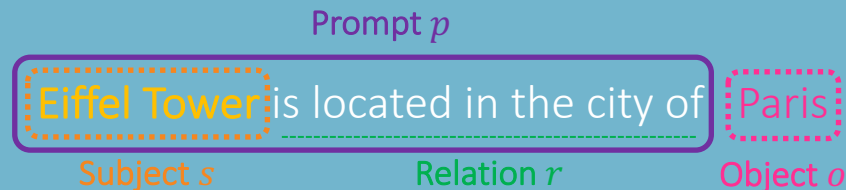
- *Factual*: concerned with **facts or contains facts**, rather than giving theories or personal interpretations.
- **Factuality**: the quality of being actual or based on facts (**“fact” to be the world knowledge**)
- **Faithfulness**: stay consistent and truthful to the provided source (**opposite to Hallucination**)
- *Factual Associations*: **causal effects** between subject and object, **based on facts (world knowledge)**
- *Factual Storage*: mechanism or some place that triggers or stores **Factual Knowledge**

Fact Representation

- *Knowledge Tuple*: $t = (s, r, o)$ where s : Subject, r : relationship, o : Object

Input and Output

- *Input*: a natural language **prompt $p = (s, r)$**
- *Output*: model's prediction of **Object o**



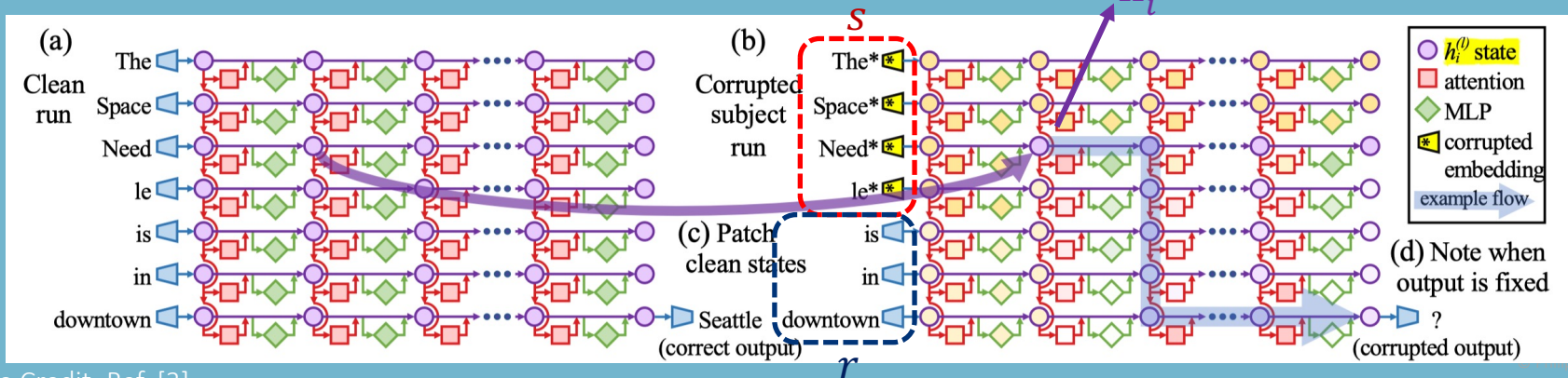
Part 1: Causal Tracing of Factual Associations

Why Causal Tracing?

- Understand Factual Associations **Eiffel Tower** is located in the city of **Paris**
- Locate the specific modules that mediate recall of a fact about a subject

How to implement Causal Tracing of Factual Associations?

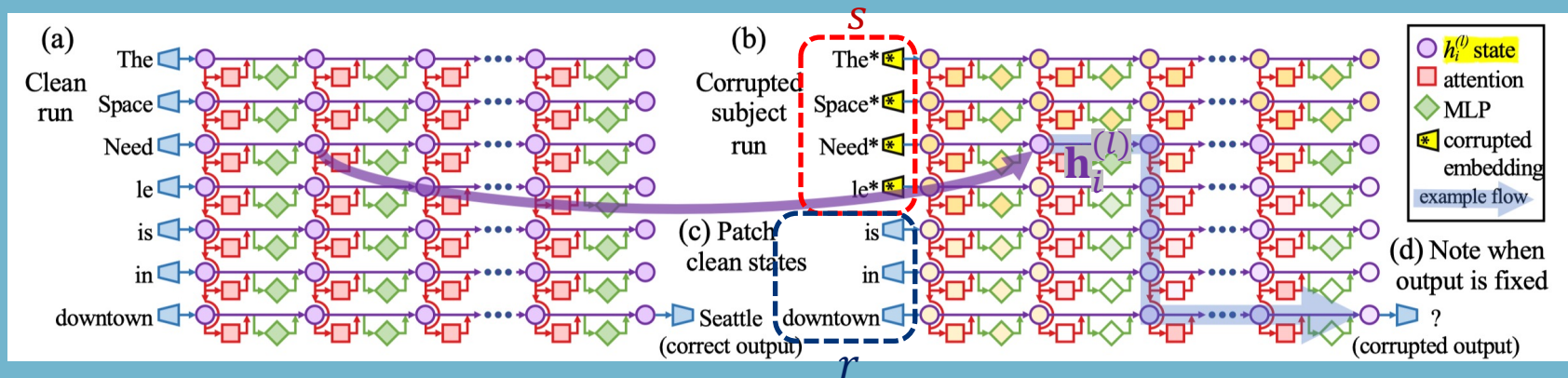
- Causal Graph and Causal Mediation Analysis



Part 1: Causal Tracing of Factual Associations



Causal Mediation Analysis: quantify the contribution of **intermediate Variables**



Clean Run $\mathbb{P}[o]$

- Clean Input $p \Rightarrow$ Hidden State $\mathbf{h}_i^{(l)}$

Corrupted Run $\mathbb{P}_*[o]$

- Noisy Input $\mathbf{h}_i^{(0)} := \mathbf{h}_i^{(0)} + \epsilon$ note: $(\epsilon \sim \mathcal{N}(0; \sigma)) \Rightarrow$ Corrupted Activations $\mathbf{h}_i^{(l)*}$

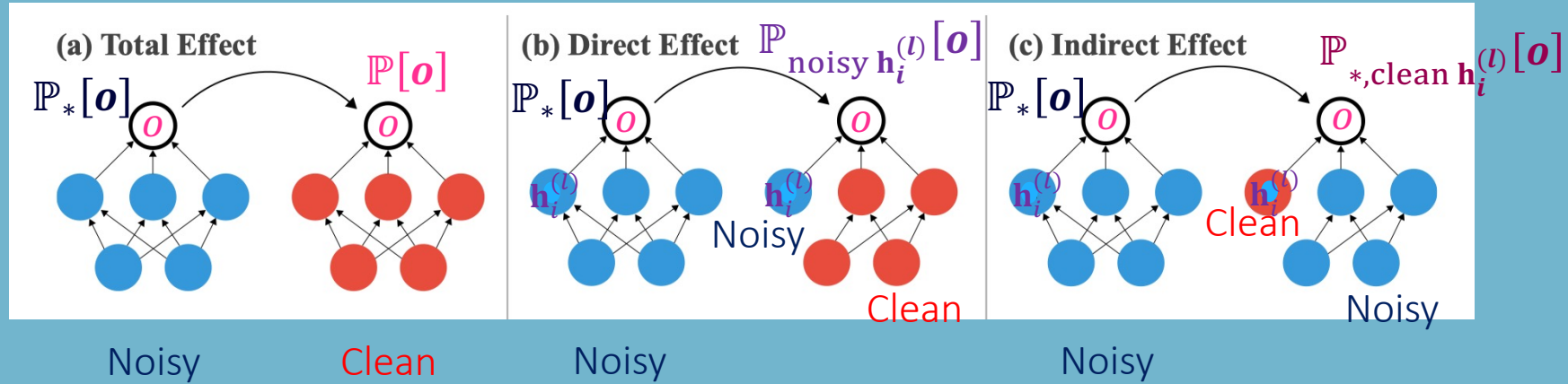
Corrupted-with-restoration Run $\mathbb{P}_{*,\text{clean } \mathbf{h}_i^{(l)}}[o]$

- Noisy Input $\mathbf{h}_i^{(0)} := \mathbf{h}_i^{(0)} + \epsilon$ except at some token \hat{l} and layer \hat{l}

Total Effect (TE)

Indirect Effect (IE)

Part 1: Causal Mediation Analysis



- Total Effect (TE) = $\mathbb{P}_*[\mathbf{o}] - \mathbb{P}[\mathbf{o}]$
 \Rightarrow change in o resulting from the intervention
- Direct Effect (DE) = $\mathbb{P}_*[\mathbf{o}] - \mathbb{P}_{\text{noisy } h_i^{(l)}}[\mathbf{o}]$
 \Rightarrow change in o resulting from performing the intervention while holding a mediator $h_i^{(l)}$ fixed
- Indirect Effect (IE) = $\mathbb{P}_*[\mathbf{o}] - \mathbb{P}_{*, \text{clean } h_i^{(l)}}[\mathbf{o}]$
 \Rightarrow change in o caused by setting $h_i^{(l)}$ to clean value, while holding others fixed

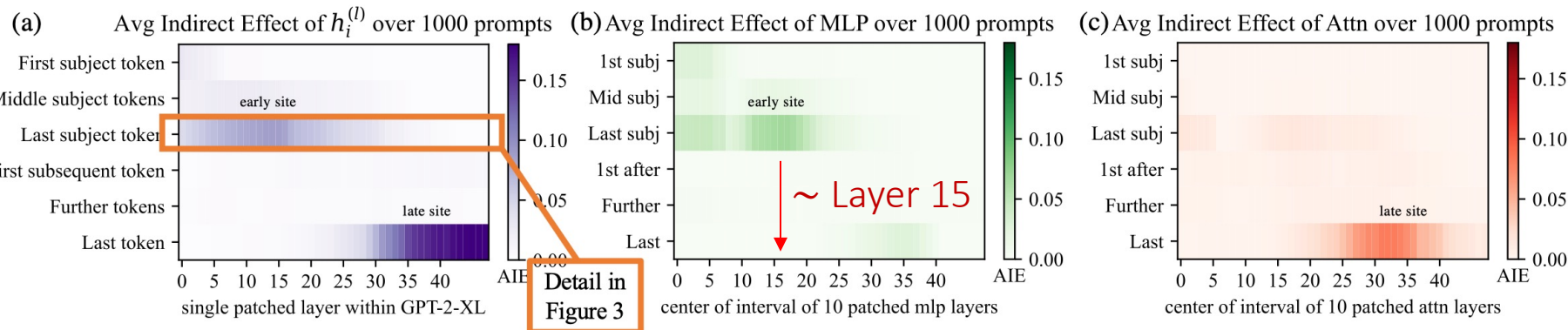
Part 1: Causal Tracing of Factual Associations



$h_i^{(L)}$

MLP

Attention



How LLMs store Factual Knowledge/Associations?

GPT-2 XL: 48 layers

MLP: contribute to the last subject token at early site and last token at late site

Attention: contribute to the last token at late site

Decisive information is accumulated across layers

Part 1: Causal Tracing of Factual Associations

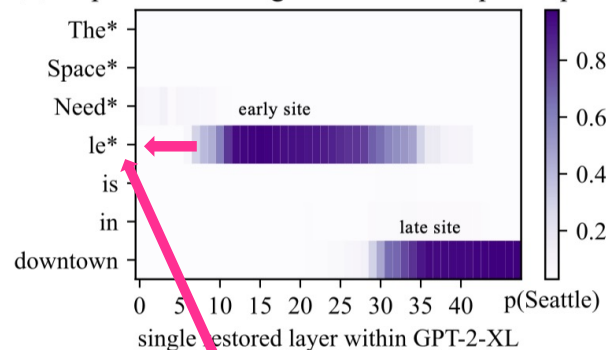


$h_i^{(l)}$

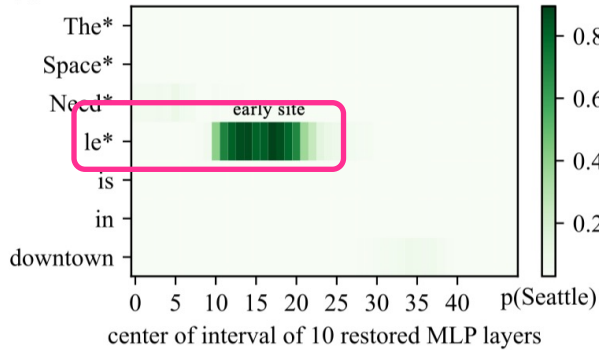
MLP

Attention

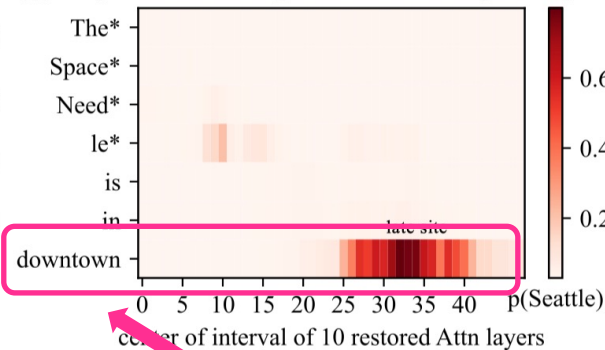
(e) Impact of restoring state after corrupted input



(f) Impact of restoring MLP after corrupted input



(g) Impact of restoring Attn after corrupted input



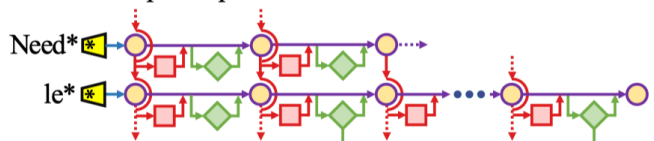
Last Subject Token

Last Token

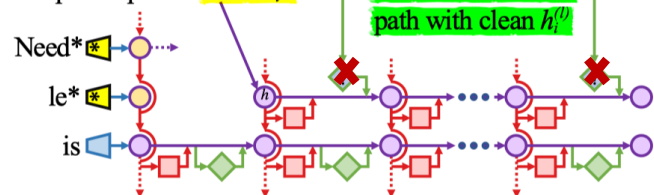
Part 1: Causal Tracing of Factual Associations

How LLMs store Factual Knowledge/Associations?

(a) baseline corrupted input condition



(b) corrupted input w/ clean $h_i^{(l)}$



(c) Causal effect of states at the early site with Attn or MLP modules severed

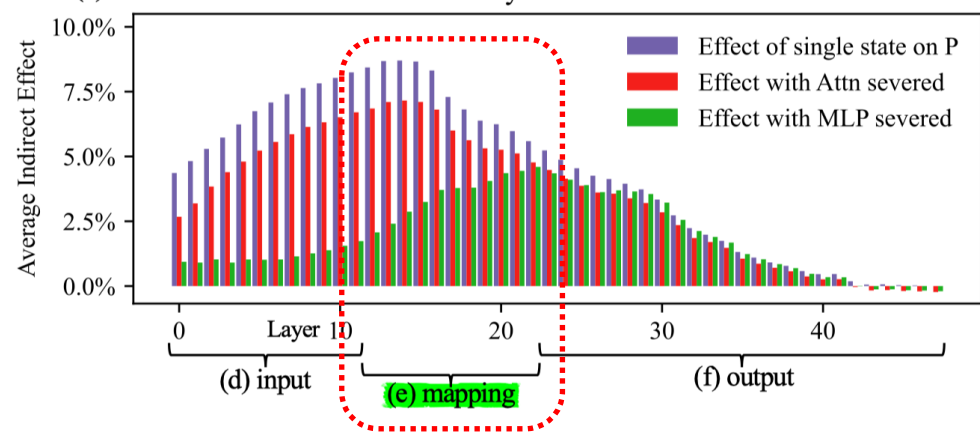
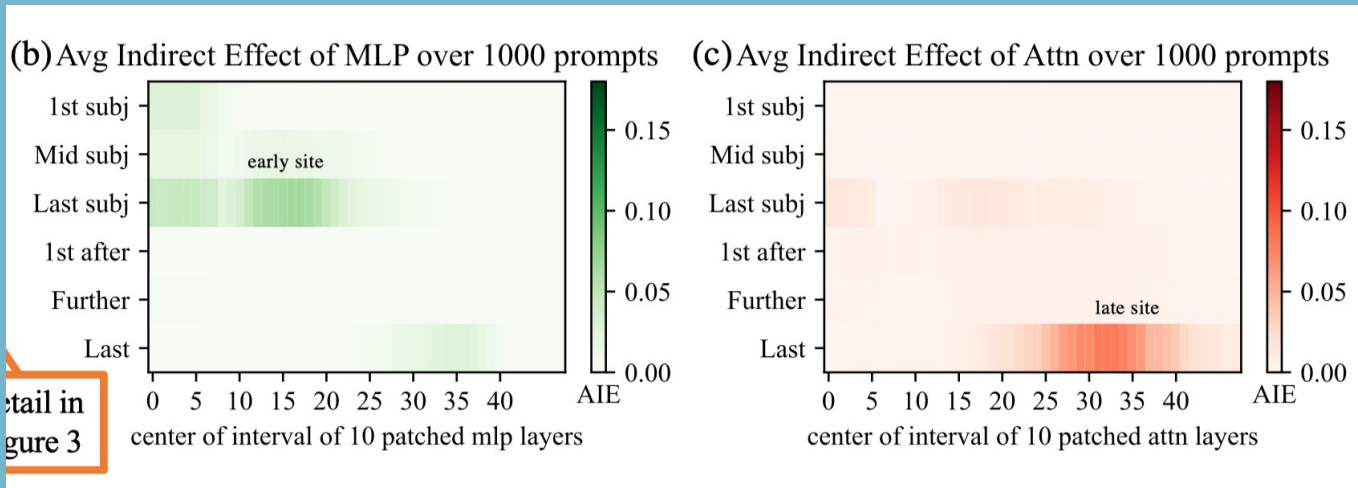


Figure 3: Causal effects with a modified computation graph. (a,b) To isolate the effects of MLP modules when measuring causal effects, the computation graph is modified. (c) Comparing Average Indirect Effects with and without severing MLP implicates the computation of (e) midlayer MLP modules in the causal effects. No similar gap is seen when attention is similarly severed.

Remove MLP or Attention \Rightarrow **MLP** module computation at **middle layers** when recalling a fact.

Part 1: Storage of Factual Associations Hypothesis



MLP Middle Layers:

- recall memorized properties about that subject
- accumulate information

Attention Layers:

- summed information is copied to the last token by attention at high layers

Part 2: Edit Weights to Understand Factual Storage



Why Edit Model Weights?

- Understand how facts are stored in weights
- Generate factual content

How to edit Model Weights?

- Rank-One Model Editing (ROME)
- By viewing $\mathbf{W}_{proj}^{(L)}$ as linear associative memory
- Update Rule:

$$\hat{\mathbf{W}} = \mathbf{W} + \Lambda (\mathbf{C}^{-1} \mathbf{K}_*)^T$$

$$\mathbf{C} = \mathbf{K} \mathbf{K}^T$$

$$\Lambda = \frac{\mathbf{V}_* - \mathbf{W} \mathbf{K}_*}{(\mathbf{C}^{-1} \mathbf{K}_*)^T \mathbf{K}_*}$$

STEP 3

Inserting the Fact

How to edit LLMs to generate Factual Recall?

$$\mathbf{W} \mathbf{K} \approx \mathbf{V} \Rightarrow \min \|\hat{\mathbf{W}} \mathbf{K} - \mathbf{V}\|$$
$$\hat{\mathbf{W}} \mathbf{K}_* = \mathbf{V}_*$$

\mathbf{K} : Key Input (e.g., Eiffel Tower)

\mathbf{V} : Value Output (e.g., Paris)

Represent the new property

(r, o^*)

\mathbf{W} : Key-Value Pair

Next: choose the appropriate \mathbf{K}_* and \mathbf{V}_*

Part 2: Edit Weights to Understand Factual Storage



STEP 1: Choose \mathbf{K}_* to represent the last subject token

- Collect Activations from a small amount of texts x that contain Subject s

$$\mathbf{K}_* = \frac{1}{N} \sum_{j=1}^N \sigma \left(\mathbf{w}_{fc}^{(l^*)} \gamma \left(\mathbf{a}_{[x_{j+s}],i}^{(l^*)} + \mathbf{h}_{[x_{j+s}],i}^{(l^*)} \right) \right)$$

STEP 2: Choose \mathbf{V}_* to recall the fact (new relation: r, o^*) $\Rightarrow \mathbf{V}_* = \operatorname{argmin}_{\mathbf{z}} (\mathcal{L}(\mathbf{z}))$

$$\mathcal{L}(\mathbf{z}) = \frac{1}{N} \sum_{j=1}^N -\log \mathbb{P}_{G(m_i^{(l^*)} := \mathbf{z})} [o^* \boxed{x_j + p}] + D_{\text{KL}} \left(\mathbb{P}_{G(m_i^{(l^*)} := \mathbf{z})} [x | p'] \parallel \mathbb{P}_G [x | p'] \right)$$

p' : “{ subject } is a”

Maximizing o^* Probability

Controlling essence drift

Current related works of Model Editing

Fine-Tuning (FT)

- applies Adam with early stopping at one layer to minimize $-\log\mathbb{P}[o^* | x]$

Constrained Fine-Tuning (FT+L)

- additionally imposes a parameter-space L_∞ norm constraint on weight changes

Knowledge Editor (KE) and MEND

- learn auxiliary models to predict weight changes

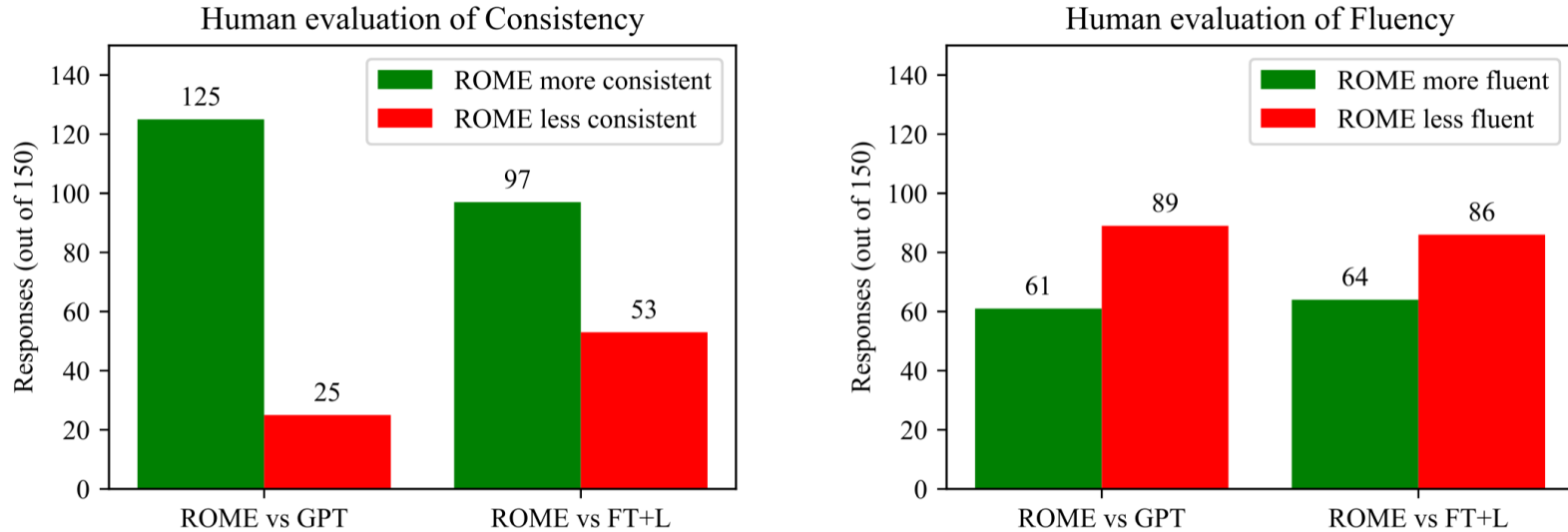


Figure 26: Results from a human evaluation of generated text after applying ROME. Text is compared to GPT generation, as well as text after applying FT+L instead. **Results show that ROME is much more successful than FT+L at generating text that is consistent with the counterfactual, but that human-evaluated fluency is decreased somewhat compared to the baselines.** Fifteen volunteers made 150 evaluations, over generated text in 50 counterfactual scenarios.

Better Factual Association Consistency but worse Generation Fluency

Potential Future Work



- Develop a **Safety Verification Method** by measuring *Factual Associations Consistency*
- Improve *Factual Associations Consistency and Generation Fluency*
- Improve *Specificity*: edited model's accuracy on an unrelated fact.

References



- [1] Ji, Ziwei, Nayeon Lee, Rita Frieske, Tiezheng Yu, Dan Su, Yan Xu, Etsuko Ishii, Ye Jin Bang, Andrea Madotto, and Pascale Fung. "[Survey of Hallucination in Natural Language Generation.](#)" ACM Computing Surveys 55, no. 12 (2023): 1-38.
- [2] Paaß, Gerhard, and Sven Giesselbach. "[Foundation Models for Natural Language Processing: Pre-trained Language Models Integrating Media.](#)" arXiv preprint arXiv:2302.08575 (2023).
- [3] Meng, Kevin, David Bau, Alex Andonian, and Yonatan Belinkov. "[Locating and Editing Factual Associations in GPT.](#)" Advances in Neural Information Processing Systems 35 (2022): 17359-17372.
- [4] Vig, Jesse, Sebastian Gehrmann, Yonatan Belinkov, Sharon Qian, Daniel Nevo, Yaron Singer, and Stuart Shieber. "[Investigating Gender Bias in Language Models using Causal Mediation Analysis.](#)" Advances in Neural Information Processing Systems 33 (2020): 12388-12401.