# Emergent Symbolic Mechanisms Support Abstract Reasoning in Large Language Models
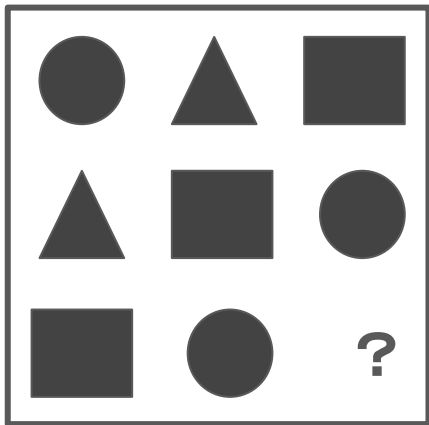
Yukang Yang[1], Declan Campbell[1], Kaixuan Huang[1], Mengdi Wang[1], Jonathan Cohen[1], Taylor Webb[2‡]

1. Princeton University  2. Microsoft Research

‡Corresponding Author: taylor.w.webb@gmail.com

https://github.com/yukang123/LLMSymbMech

# Introduction



Raven's Progressive Matrix

Abstract Reasoning

- Large Language Models (LLMs) have shown impressive performance on various human reasoning tasks including abstract (analogical) reasoning[*].

- Some studies questioned the robustness of LLMs' reasoning abilities[^].

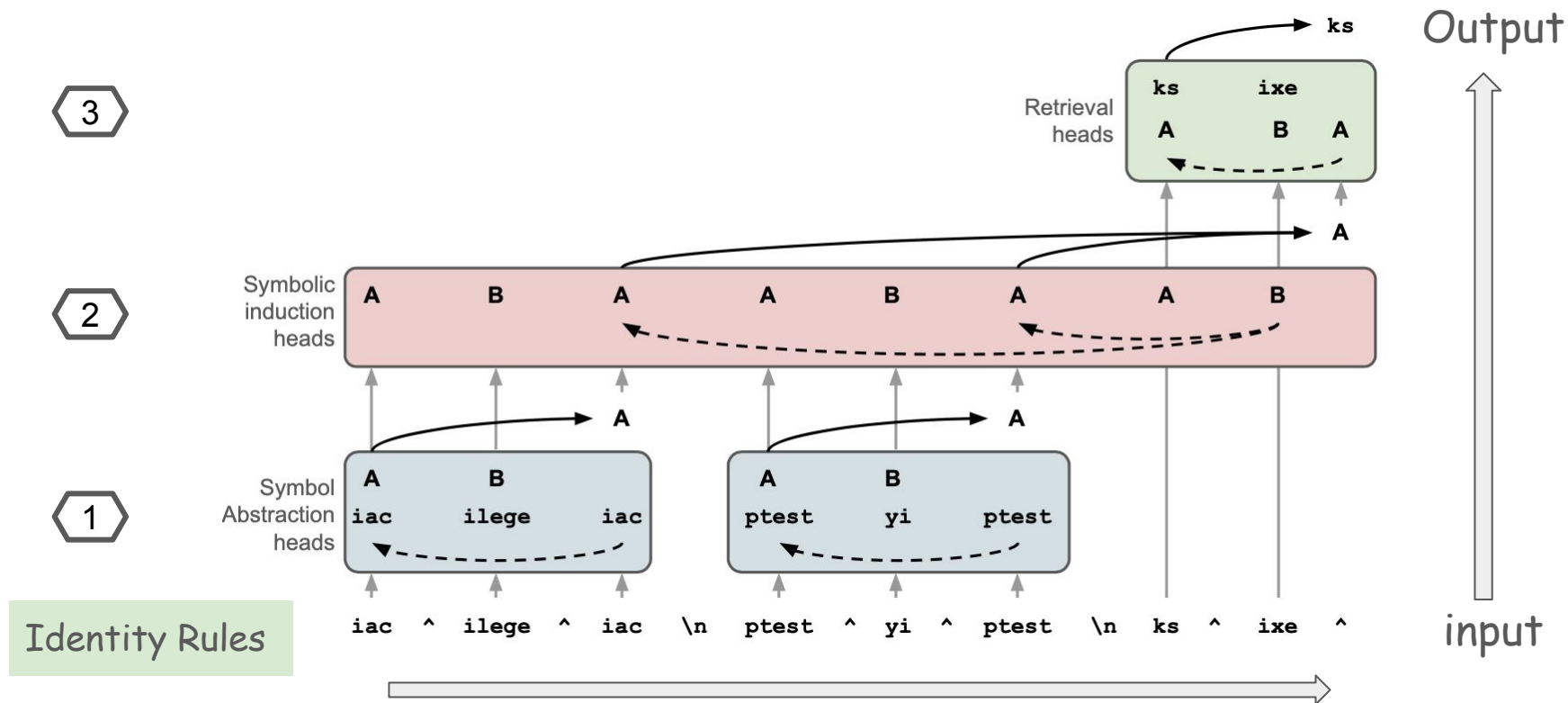structured human-like **symbolic processing** Mechanism?
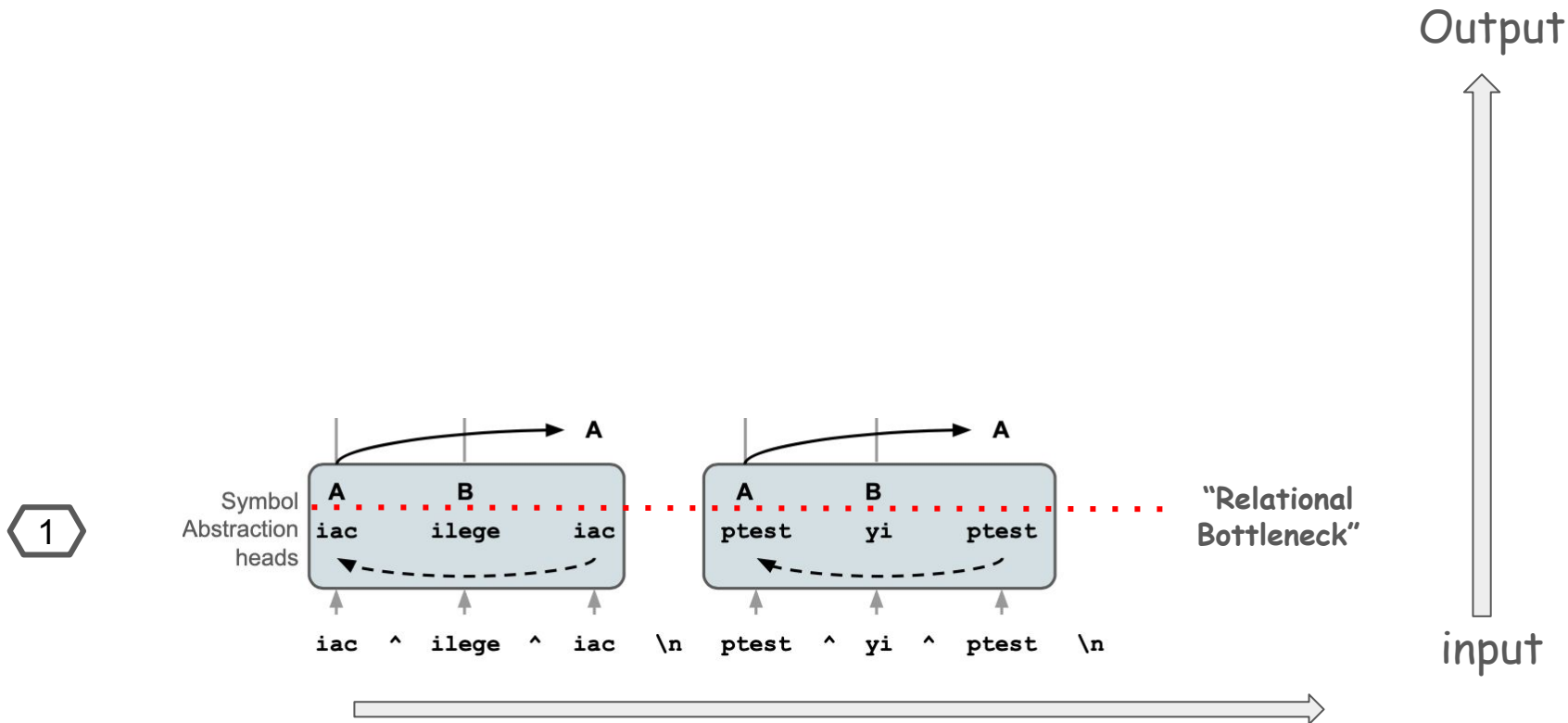
OR

Statistically approximating training data?

*Taylor, Webb, et al,. "Emergent analogical reasoning in large language models." *Nature Human Behaviour* 7.9 (2023): 1526-1541.
^Lewis, M. and Mitchell, M. Evaluating the robustness of analogical reasoning in large language models. arXiv:2411.14215, 2024.

# Our Findings : **Emergent** Symbolic Processing Mechanism in LLMs

# **Emergent** Symbolic Processing Mechanism in LLMs

Output

⬡ 1

Symbol Abstraction heads

| A | B | A |
|---|---|---|
| iac | ilege | iac |

| A | B | A |
|---|---|---|
| ptest | yi | ptest |

"Relational Bottleneck"

iac  ^  ilege  ^  iac  \n  ptest  ^  yi  ^  ptest  \n

input

# **Emergent** Symbolic Processing Mechanism in LLMs

Output

input

Operating in the
**symbol** space

Symbolic
induction
heads

② 

A   B   A   A   B   A   A   B   A

Symbol
Abstraction
heads

① 

A       B       A
iac   ilege   iac

A       B       A
ptest   yi   ptest

iac  ^  ilege  ^  iac  \n  ptest  ^  yi  ^  ptest  \n  ks  ^  ixe  ^

# Emergent Symbolic Processing Mechanism in LLMs

**Our Analyses:** (please refer to the paper)

1. Causal Mediation Analyses (CMA)

2. Attention Analyses

3. Representation Similarity Analyses (RSA)

4. Ablation Studies

5. Comparison with Induction Heads and Function Vectors

**Our Analyses:** (please refer to the paper)

1. **Causal Mediation Analyses (CMA)**

2. Attention Analyses

3. Representation Similarity Analyses (RSA)

4. Ablation Studies

5. Comparison with Induction Heads and Function Vectors

# Causal Mediation Analyses (CMA)

1. Design a context pair ($c_1$, $c_2$) to isolate either <u>abstract symbols</u> or <u>literal tokens</u>

**Abstract Context Pair** : test whether the embedding represents **abstract symbols**

(used for *symbol abstraction heads* and *symbolic induction heads*)

$$c_1^{abstract} = La\char`^Li\char`^La\backslash nTe\char`^To\char`^Te\backslash nHi\char`^Ha\char`^$$

$$c_2^{abstract} = Li\char`^La\char`^La\backslash nTo\char`^Te\char`^Te\backslash nHa\char`^Hi\char`^$$

Answer

$y_{c_1}$   $Hi$

$y_{c_2}$   $Hi$

Same Token, Different Symbols/Rules (A vs B)

# Causal Mediation Analyses (CMA)

1. Design a context pair $(c_1, c_2)$ to isolate either <u>abstract symbols</u> or <u>literal tokens</u>

**Token Context Pair** : test whether the embedding represents **literal tokens**

(used for *Retrieval heads*)

Answer

$$c_1^{token} = La\hat{}Li\hat{}La\backslash nTe\hat{}To\hat{}Te\backslash nHi\hat{}Ha\hat{}$$
$$c_2^{token} = La\hat{}Li\hat{}La\backslash nTe\hat{}To\hat{}Te\backslash nHa\hat{}Hi\hat{}$$
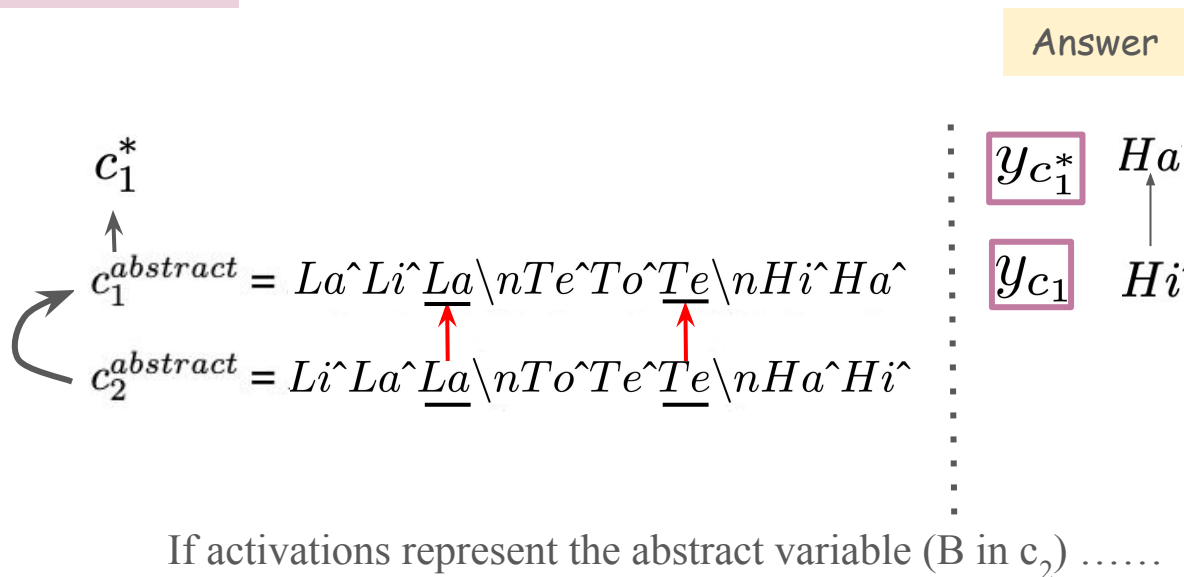
$y_{c_1}$  $Hi$

$y_{c_2}$  $Ha$

Same Rule, Different Tokens

# Causal Mediation Analyses

2. Activation Patching at <u>certain token positions</u>:

Replace attention head outputs in context $c_1$ with the corresponding activations from $c_2$

Answer

$$c_1^*$$

$$c_1^{abstract} = La\hat{\ }Li\hat{\ }\underline{La}\backslash nTe\hat{\ }To\hat{\ }\underline{Te}\backslash nHi\hat{\ }Ha\hat{\ }$$

$$c_2^{abstract} = Li\hat{\ }La\hat{\ }\underline{La}\backslash nTo\hat{\ }Te\hat{\ }\underline{Te}\backslash nHa\hat{\ }Hi\hat{\ }$$

$y_{c_1^*}$ $Ha$

$y_{c_1}$ $Hi$

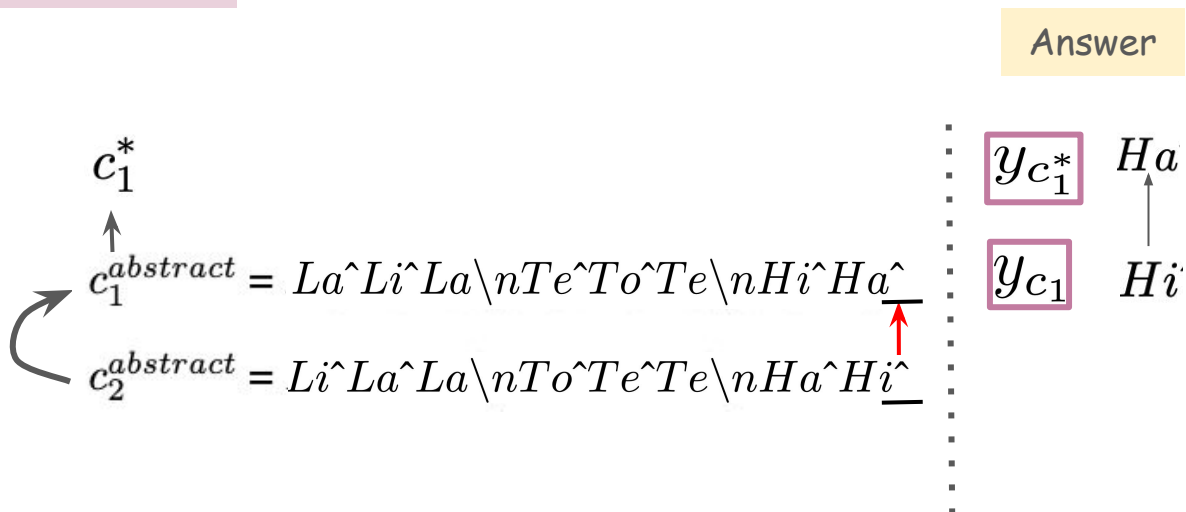If activations represent the abstract variable (B in $c_2$) ……

# Causal Mediation Analyses

2. Activation Patching at <u>certain token positions</u>:

   Replace attention head outputs in context $c_1$ with the corresponding activations from $c_2$

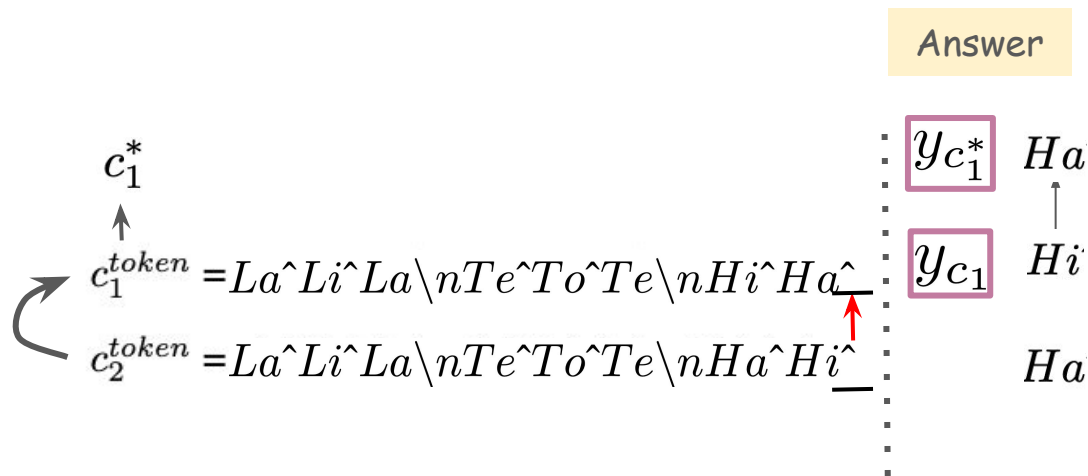<span style="background-color:#e6c8d8">symbolic induction heads</span>

<span style="background-color:#fdf0c8">Answer</span>

$$c_1^*$$

$$c_1^{abstract} = La\hat{}Li\hat{}La\backslash nTe\hat{}To\hat{}Te\backslash nHi\hat{}Ha\hat{}\underline{\phantom{x}}$$

$$c_2^{abstract} = Li\hat{}La\hat{}La\backslash nTo\hat{}Te\hat{}Te\backslash nHa\hat{}Hi\hat{}\underline{\phantom{x}}$$

$$\boxed{y_{c_1^*}} \quad Ha$$

$$\boxed{y_{c_1}} \quad Hi$$

If activations represent the abstract variable (B in $c_2$) ……

# Causal Mediation Analyses

2. Activation Patching at <u>certain token positions</u>:

   Replace attention head outputs in context $c_1$ with the corresponding activations from $c_2$

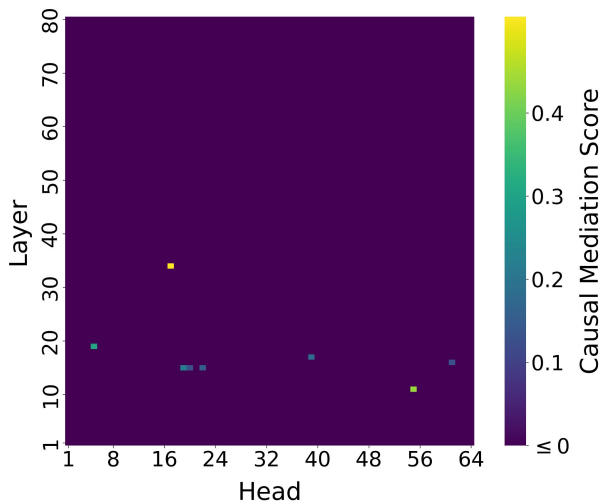<div style="background-color:#e8d5e0; display:inline-block; padding:4px 12px;">Retrieval heads</div>

Answer

$$c_1^*$$

$$\uparrow$$

$$c_1^{token} = La\char`\^Li\char`\^La\backslash nTe\char`\^To\char`\^Te\backslash nHi\char`\^Ha\char`\^\_$$

$$c_2^{token} = La\char`\^Li\char`\^La\backslash nTe\char`\^To\char`\^Te\backslash nHa\char`\^Hi\char`\^\_$$

$\boxed{y_{c_1^*}} \quad Ha$

$\boxed{y_{c_1}} \quad Hi$

$\quad Ha$

If activations represent the literal token for the answer ($Ha$ in $c_2$) ……
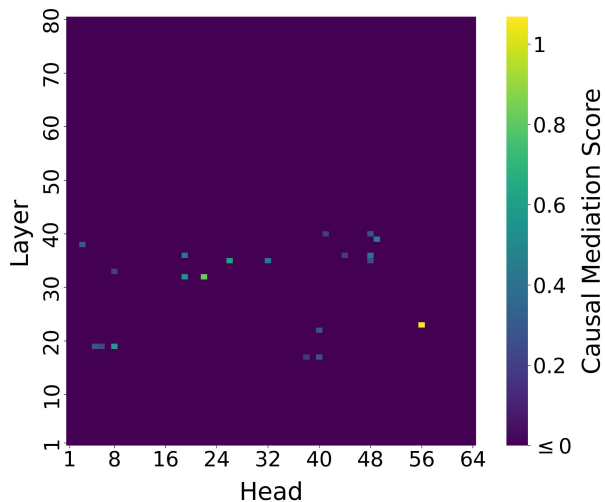
# Causal Mediation Analyses

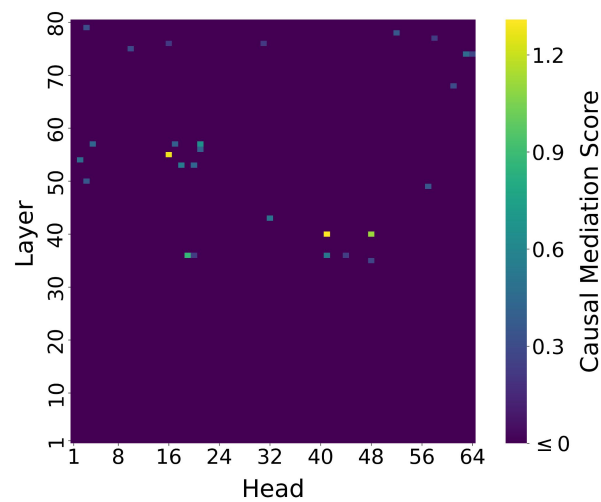3. Measure whether the causal effects on $c_1$ comply with the hypotheses

Llama-3.1 70B



| Symbol Abstraction Head | Symbolic Induction Head | Retrieval Head |

- Defined a score based on the changes in output logits for the answers as a measure of causal effects.
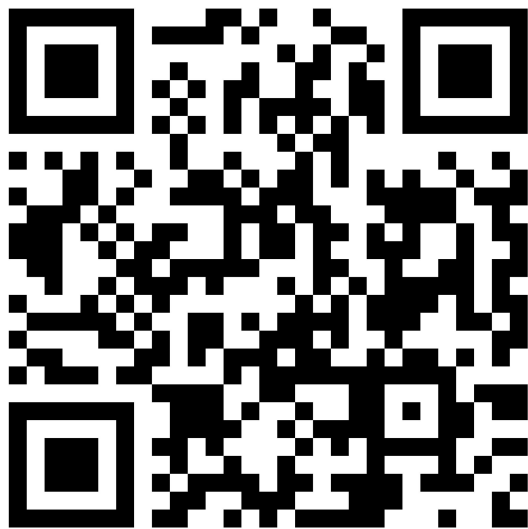- Conducted multi-hypothesis permutation tests to select significant ones

**Our Analyses:** (please refer to the paper)

❖ Experiments with Llama-3.1 70B on identity rule tasks

1. Causal Mediation Analyses (CMA)

2. Attention Analyses

3. Representation Similarity Analyses (RSA)

4. Ablation Studies

5. Comparison with Induction Heads and Function Vectors

**Our Analyses:** (please refer to the paper)

❖ Evaluating **more LLMs** on identity rule tasks

➤ Tested *13* models across *4* model families (GPT-2, Gemma-2, Qwen2.5, Llama-3.1)

➤ Found similar robust symbolic mechanisms in *3* model families with GPT-2 as an exception

➤ GPT-2 models showed low generation accuracy while symbol abstraction heads rarely emerged.

❖ Studying **more** complex abstract **reasoning tasks**

➤ *Letter string analogies* and *verbal analogies*

➤ Identified similar three-stage symbolic processing structures through CMA.

➤ Different tasks may involve different attention heads to implement symbol processing.

Please Check out our paper!