

Truth Neurons

Haohang Li Yupeng Cao Yangyang Yu Jordan W. Suchow* Zining Zhu*

Stevens Institute of Technology

{hli1113, ycao33, yyu44, jws, zzhu41}@stevens.edu

Abstract

Despite their remarkable success and deployment across diverse workflows, language models sometimes produce untruthful responses. Our limited understanding of how truthfulness is mechanistically encoded within these models jeopardizes their reliability and safety. In this paper, we propose a method for identifying representations of truthfulness at the neuron level. We show that language models contain *truth neurons*, which encode truthfulness in a subject-agnostic manner. Experiments conducted across models of varying scales validate the existence of truth neurons, confirming that the encoding of truthfulness at the neuron level is a property shared by many language models. The distribution patterns of truth neurons over layers align with prior findings on the geometry of truthfulness. Selectively suppressing the activations of truth neurons found through the TruthfulQA dataset degrades performance both on TruthfulQA and on other benchmarks, showing that the truthfulness mechanisms are not tied to a specific dataset. Our results offer novel insights into the mechanisms underlying truthfulness in language models and highlight potential directions toward improving their trustworthiness and reliability.

1 Introduction

Language models have demonstrated remarkable text-generation capabilities across various tasks [22, 28, 53], but they struggle to consistently produce correct outputs in certain question-answering scenarios [17, 18]. The struggle arises partly because language models lack sufficient relevant knowledge about specific questions in their pretrained data [6]. Moreover, language models may generate incorrect answers despite recognizing the incorrectness of the responses [56]. For instance, prior research has shown that language models aligned with human feedback tend to accommodate users’ incorrect responses, even when the models initially identify these responses as false [48]. Although the correctness of language models can be substantially improved through self-consistency checking [31], post-training [44], and optimizing decode strategies [8, 7], it is still unknown whether there exists a *truth mechanism*, a special mechanism within language models that drives the generation of accurate answers.

Research on mechanistic interpretability has begun to probe representations of truthfulness through analyses of hidden states: Orgad et al. [35] applied linear probes to reveal meaningful patterns of truth-related encoding. Marks and Tegmark [32] identified specific tokens and layers involved in truthfulness and demonstrated a linear encoding of truth and falsehood using principal component analysis (PCA). Ferrando et al. [12] used sparse autoencoders (SAEs) to identify the features related to entity awareness and hallucination.

Despite these advancements, neuron-level mechanisms of truthfulness remain unknown. The neuron is a fundamental level of analysis in both the human brain and Transformer-based neural networks.

*Co-corresponding authors.

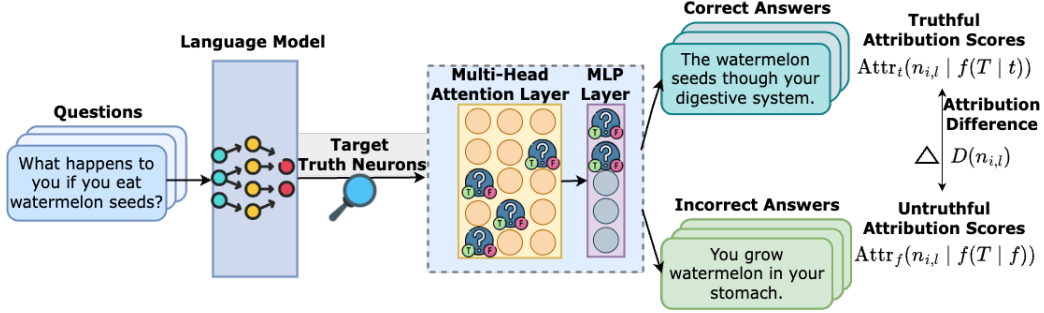


Figure 1: Overview of our method that detects the truth neurons.

For example, specific neurons in the human brain (e.g., those in the dorsolateral and ventrolateral prefrontal cortex) selectively activate when performing certain cognitive operations, such as evaluating the truthfulness of particular events [20, 37, 19, 21]. Analogous to these observations in the human brain, the transformer-based models that underlie language models also exhibit functional specialization. Transformers are believed to activate distinct regions selectively, facilitating interactions necessary for informed decision-making, such as true-or-false judgments [27, 49, 24, 41, 25]. At the neuron level, recent research has also observed the knowledge storage and retrieval mechanisms to varying extents [9, 33], but as we will show, the mechanisms to process truth differ from those of the knowledge entities. Truth mechanisms are not localized to specific entities (even datasets), whereas the knowledge storage is localized to each data entry.

Here, we develop a method informed by neuroscience and interpretability research to detect *truth neurons*, specialized truth-processing structures within language models. Our method starts with an axiomatic attribution [42], using integrated gradients to measure neuron attribution scores for truthful vs. untruthful responses. We identify candidate neurons positively contributing to truthfulness and negatively correlated with untruthfulness. We then apply a systematic filtering procedure to select a small subset of neurons causally linked to truthfulness representations. Upon suppressing these identified *truth neurons*, we observe a statistically significant reduction in accuracy on the TruthfulQA benchmark [30]. Further analysis reveals that this reduction is not biased toward any specific category, suggesting that these neurons encode a general, category-agnostic representation of truthfulness. Additionally, we demonstrate that the influence of the truth neurons generalizes effectively to other truthfulness benchmarks. As we will show in our experiments, identifying and analyzing truthfulness mechanisms at neuronal granularity reveals insights that deepen our understanding of truthfulness representations in Transformer-based language models.

In summary, our work makes the following contributions:

- We propose a novel method to identify truth neurons. By analyzing neuron attributions, we successfully isolate a small subset of neurons whose activations have a statistically significant impact on the model’s ability to discern truthfulness (Section 2).
- Through carefully designed experiments, we demonstrate that the identified neurons encode general, example-agnostic representations of truthfulness, and that their influence generalizes effectively to other out-of-distribution truthfulness benchmarks (Section 3.3 & Section 3.4).
- Finally, we investigate the distribution patterns of these identified truth neurons across model layers, observing a consistent pattern that aligns closely with existing findings (Section 3.5).

We believe our results offer insights into improving the trustworthiness and safe deployment of language models, highlighting promising future directions for enhancing model alignment with truthfulness.

2 Methodology

In this section, we propose an integrated gradient-based approach [42] to systematically identify and isolate neurons causally associated with a model’s ability to discern the truthfulness of factual statements. Within Transformer architectures, feed-forward (MLP) layers have been characterized

as key-value memory structures closely tied to factual knowledge recall [13]; attention heads have similarly been linked to truthfulness representations [29]. Therefore, we extend neuron attribution analyses to encompass both MLP and attention modules across all intermediate layers.

2.1 Preliminaries

As our goal is to identify neurons correlated with the model’s truthfulness behavior, integrated gradient is a suitable tool, as it satisfies desirable axioms and effectively quantifies each neuron’s contribution to model behavior. Following the setup of Sundararajan et al. [42]. Let $\mathbf{X} \in \mathbb{R}^n$ be the input tensor of the neural network, $\mathbf{X}' \in \mathbb{R}^n$ be the baseline input tensor required by the method, and $\mathbf{f} : \mathbb{R}^n \rightarrow \mathbb{R}$ denote the function representing the neural network. Additionally, define $n_{i,l}^{\text{input}}$ as the intermediate neuron activation output at layer l and index i , with $n_{i,l}^{\text{baseline}}$ representing the corresponding activation when the baseline input is applied. The integrated gradient method computes neuron attribution as a path integral along the straight-line path $\gamma(\alpha)$ from \mathbf{X}' to \mathbf{X} , where α represents the incremental interpolation parameter indicating progress along the path:

$$\gamma(\alpha) = n_{i,l}^{\text{baseline}} + \alpha(n_{i,l}^{\text{input}} - n_{i,l}^{\text{baseline}}), \alpha \in [0, 1] \quad (1)$$

$$\text{Attr}(n_{i,l} \mid \mathbf{f}(\mathbf{X})) := \int_0^1 \frac{\partial \mathbf{f}(\gamma(\alpha))}{\partial \gamma(\alpha)} d\alpha = \int_0^1 \frac{\partial \mathbf{f}(\gamma(\alpha))}{\partial \gamma(\alpha)} \times \frac{d\gamma(\alpha)}{d\alpha} d\alpha \quad (2)$$

$$= (n_{i,l}^{\text{input}} - n_{i,l}^{\text{baseline}}) \int_0^1 \frac{\partial \mathbf{f} \left(n_{i,l}^{\text{baseline}} + \alpha (n_{i,l}^{\text{input}} - n_{i,l}^{\text{baseline}}) \right)}{\partial n_{i,l}} d\alpha. \quad (3)$$

Intuitively, integrated gradient attribution quantifies a neuron’s contribution to the final prediction by measuring how the predicted probability changes as the neuron’s activation is gradually shifted from its baseline value toward its activation in the actual input. In the computation, the integral is approximated by a Riemann sum:

$$\text{Attr}(n_{i,l} \mid \mathbf{f}(\mathbf{X}))^{\text{Approx}} = \frac{n_{i,l}^{\text{input}} - n_{i,l}^{\text{baseline}}}{m} \sum_{k=1}^m \frac{\partial \mathbf{f} \left(n_{i,l}^{\text{baseline}} + \frac{k}{m} (n_{i,l}^{\text{input}} - n_{i,l}^{\text{baseline}}) \right)}{\partial n_{i,l}}, \quad (4)$$

where m is the step parameter that controls the approximation precision.

2.2 Identifying Truth Neurons

Notation. For each question q , the dataset provides one correct answer t and one incorrect answer f , with the incorrect answer closely matching the length and format of the correct answer whenever possible. We construct the input prompt \mathbf{T} by appending these two answers after the question in randomized order, labeling them as options A and B. Additionally, an instruction i explicitly prompts the model to select the option that correctly answers the question. We can then denote a dataset \mathcal{D} with N questions as:

$$\mathcal{D} = \{\mathbf{T}^{(k)}\}_{k=1}^N = \{\langle q, t, f, i \rangle^{(k)}\}_{k=1}^N, \quad (5)$$

where k indexes the dataset \mathcal{D} .

Accounting for upper and lower cases. Let \mathcal{M} denote the language model’s output probability distribution. We observed that language models frequently interchange the uppercase and lowercase forms of output labels. To cover both cases, we define the prediction probability \mathbf{f} as the sum of both the uppercase probability and the lowercase probability. For example, when the correct answer is labeled A, the probability for the correct response is:

$$\mathbf{f}(\mathbf{T} \mid t) = \mathcal{M}(\hat{y} = \text{A} \mid \mathbf{T}) + \mathcal{M}(\hat{y} = \text{a} \mid \mathbf{T}). \quad (6)$$

Similarly, the probability for the incorrect response labeled B is:

$$\mathbf{f}(\mathbf{T} \mid f) = \mathcal{M}(\hat{y} = \mathbf{B} \mid \mathbf{T}) + \mathcal{M}(\hat{y} = \mathbf{b} \mid \mathbf{T}). \quad (7)$$

This definition applies analogously when the correct answer is labeled B and the incorrect answer A. Note that for both the correct and the incorrect answers, we query the probabilities from the same distribution (i.e., the same prompt \mathbf{T}), avoiding the lexical biases.

Deconfounding untruthfulness. For a given neuron $n_{i,l}$ at the i th position and the l th layer, applying integrated gradients to the input with respect to the correct and incorrect responses yields $\text{Attr}_t(n_{i,l} \mid \mathbf{f}(\mathbf{T} \mid t))$ and $\text{Attr}_f(n_{i,l} \mid \mathbf{f}(\mathbf{T} \mid f))$, the two corresponding attribution scores.

We denote by $\text{Attr}_t^{\text{Avg}}$ and $\text{Attr}_f^{\text{Avg}}$ the average truthful and untruthful attribution scores computed over N examples, respectively. We further define the attribution difference for a single example as $D(n_{i,l})$, and the average attribution difference across all examples in the dataset as $\bar{D}(n_{i,l})$:

$$D(n_{i,l}) = \text{Attr}_t(n_{i,l} \mid \mathbf{f}(\mathbf{T} \mid t)) - \text{Attr}_f(n_{i,l} \mid \mathbf{f}(\mathbf{T} \mid f)) \quad (8)$$

$$\bar{D}(n_{i,l}) = \frac{1}{N} \sum_{j=1}^N (\text{Attr}_t(n_{i,l} \mid \mathbf{f}(\mathbf{T} \mid t)) - \text{Attr}_f(n_{i,l} \mid \mathbf{f}(\mathbf{T} \mid f))) = \text{Attr}_t^{\text{Avg}} - \text{Attr}_f^{\text{Avg}} \quad (9)$$

Empirically, the signs of the truthful attribution scores $\text{Attr}_t(n_{i,l} \mid \mathbf{f}(\mathbf{T} \mid t))$ and the untruthful attribution scores $\text{Attr}_f(n_{i,l} \mid \mathbf{f}(\mathbf{T} \mid f))$ can be categorized into four distinct scenarios, which consequently determine the signs of $D(n_{i,l})$:

- [1] **Both positive:** The neuron positively contributes to both correct and incorrect responses; the overall attribution difference depends on the relative magnitudes of these contributions.
- [2] **Truthful negative, untruthful positive:** The neuron predominantly supports untruthful responses, negatively contributing to truthfulness and positively correlating with untruthfulness. This combination results in a negative attribution difference, indicating it is not a truth neuron.
- [3] **Truthful positive, untruthful negative:** The neuron supports truthfulness, positively contributing to truthful responses and negatively correlating with untruthful responses. This combination yields a large positive attribution difference, clearly indicating a truth neuron.
- [4] **Both negative:** The neuron negatively contributes to both responses; the attribution difference depends on the strength of each negative contribution.

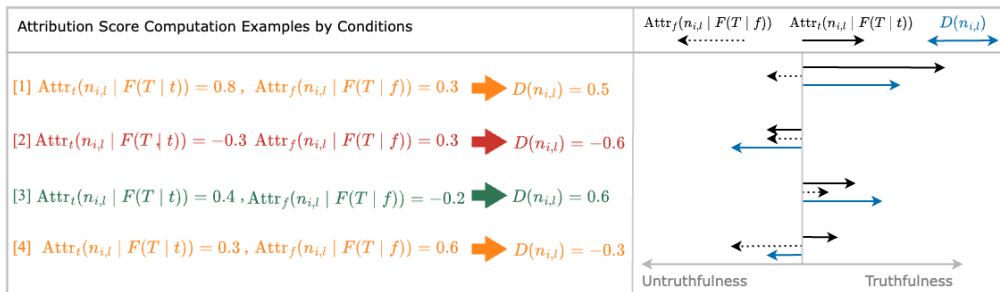


Figure 2: Examples of attribution score computation. The left side shows example attribution scores for truthful and untruthful responses, and the right side shows the resulting attribution differences. Colors correspond to the four scenarios discussed above. In case [1], competing attributions result in a positive difference, indicating a positive correlation with truthfulness, while case [4] illustrates the opposite situation. Case [2] indicates a clear bias toward untruthfulness, whereas case [3] shows a strong bias toward truthfulness.

Hypothesis testing against randomness. To test whether a neuron consistently encodes truthfulness-related information, we conducted a Student’s t -test for $\bar{D}(n_{i,l})$ against 0. The null and alternative hypotheses are defined below. If truthfulness-related information is successfully encoded, the null hypothesis will be rejected, and the alternative hypothesis will be accepted; otherwise, the reverse

will hold.

$$\begin{aligned} H_0 : \bar{D}(n_{i,l}) &\approx \text{Attr}_t^{\text{Avg}} - \text{Attr}_f^{\text{Avg}} + \epsilon = 0, \\ H_a : \bar{D}(n_{i,l}) &> 0. \end{aligned} \tag{10}$$

where ϵ is assumed to resemble random noise due to averaging over diverse inputs varying in semantics and syntax, likely activating different neurons. We applied the Bonferroni correction to the t -tests to mitigate the inflation of Type I errors caused by the multiple comparisons problem.

2.3 Systematic Filtering for Dataset and Attributions

To more accurately and efficiently identify the truth neurons of interest, we applied additional filtering steps to both the dataset and the neuron activations.

Manipulation check. We conducted a manipulation check to ensure we were probing neurons that accurately reflect the truthfulness of the language model. Specifically, we retained only those examples for which the model can answer correctly. If the model fails to correctly distinguish between truthful and untruthful responses, it indicates a lack of the necessary knowledge regarding truthfulness. Consequently, any neuron-level probing in such cases would not yield meaningful insights into the underlying mechanism of truthfulness.

Systematic filtering. To efficiently identify a candidate set of truth neurons, we follow a refining approach similar to that described in Dai et al. [9]. Specifically, we consider only those neurons whose attribution differences $D(n_{i,l})$ are notably salient across the examples. The filtering process involves two main steps. First, for each example and each layer type, we identify the maximum neuron activations across all layers and retain only those whose activations exceed an adaptive threshold set at $t\%$ of this maximum activation. Second, after identifying the most salient neurons per example for each layer type, we further require that neurons consistently remain among the most salient across at least $p\%$ of examples—referred to as the share threshold. This ensures that the selected neurons reliably represent truthfulness that generalizes across examples rather than being tied to specific input features or triggered by sporadic activations.

Adjustment to avoid double-dipping. Threshold-based neuron identification methods may suffer from non-independence errors due to the reuse of the same dataset for both neuron selection and subsequent statistical analyses, a problem known as “double-dipping” or circular analysis in statistics [26]. To avoid double-dipping, we adopted a strategy recommended by Vul et al. [47]: we split the dataset into two halves, using the first half to select the neurons and the second half to conduct statistical tests. In this way, the selection and statistical analysis procedures are separate.

3 Experiments and Results

To verify the existence of truth neurons and determine whether they faithfully represent truthfulness, we propose the following three research questions (RQs):

- **RQ1:** Do truth neurons exist across language models?
- **RQ2:** Do truth neurons identified using TruthfulQA generalize beyond that dataset?
- **RQ3:** What is the distribution pattern over layers for truth neurons within language models?

3.1 Experiment Setup

We conduct experiments using six state-of-the-art open-source models across various parameter scales to demonstrate the generalizability and robustness of our method. Specifically, we include Llama-3.2-3B-Instruct [14] and Qwen-2.5-3B-Instruct [38] as representatives of small-scale models; Llama-3.1-8B-Instruct and OLMo-2-7B-Instruct [34] as medium-scale models; and Mistral-Nemo-Instruct [43] and OLMo-2-13B-Instruct as examples of relatively large-scale models. To ensure fairness, we employ a consistent, standardized instruction prompt across all models for truth neuron identification, detailed in Figure 9. The integrated gradient method is approximated using $m = 20$ interpolation steps, and the share threshold is set to $p = 40\%$. Since attribution scales vary across models, the adaptive threshold ($t\%$) requires manual tuning. We observed that excessively high thresholds filter out too many neurons, resulting in minimal or negligible performance impacts upon

Model	Baseline	Suppressed Random Neurons	Suppressed Truthful Neurons	
	Acc. (%)	Acc. (%)	Acc. (%)	# of Neurons
Qwen2.5-3B-Instruct	65.67 \pm 0.67	65.91 \pm 1.08	58.59 \pm 0.68*	35
Llama-3.2-3B-Instruct	55.55 \pm 0.76	55.47 \pm 1.49	49.90 \pm 1.00*	114
OLMo-2-1124-7B-Instruct	50.76 \pm 0.91	51.42 \pm 0.96	49.38 \pm 1.12*	655
Llama-3.1-8B-Instruct	62.15 \pm 0.94	62.10 \pm 1.15	43.31 \pm 0.88*	37
Mistral-Nemo-Instruct-2407 (12B)	58.06 \pm 0.99	58.46 \pm 1.08	50.04 \pm 1.25*	181
OLMo-2-1124-13B-Instruct	61.89 \pm 0.63	61.85 \pm 0.71	49.35 \pm 1.16*	75

Table 1: Number of truth neurons identified under the specified hyperparameter setup, along with accuracy (Acc.) comparisons among the baseline, random-neuron suppression, and truth-neuron suppression conditions. Bold values marked with * indicate statistically significant accuracy reductions ($p < 0.05$) from the baseline to the truth-neuron suppression condition across 10 repetitions. Accuracy is reported in percentage (%).

suppression. Conversely, thresholds set too low include numerous neurons that may be unrelated to truthfulness, whose suppression significantly impairs the model’s instruction-following abilities and hinders accurate evaluation. The criteria guiding threshold selection and specific hyperparameter values for each model are provided in [Section A.3](#). The experiments are conducted with 4xNVIDIA H100 and 1xNVIDIA H200.

3.2 Datasets

TruthfulQA: To identify truthfulness representations at the neuronal level, we use the TruthfulQA dataset introduced by Lin et al. [30]. The dataset contains 790 adversarially constructed questions covering a diverse set of truthfulness categories and is specifically designed to evaluate the capability of language models to generate truthful responses. We use the updated binary-choice evaluation framework following the details outlined in [Appendix A.2](#).

TriviaQA and MMLU: To verify whether neurons identified using TruthfulQA generalize as faithful representations of truthfulness, we evaluate performance on two additional datasets employed to measure the truthfulness [29, 51, 2]: TriviaQA [23] and MMLU [15]. TriviaQA is a question-answering dataset spanning diverse topics, while MMLU is a benchmark assessing a language model’s factual knowledge across 57 subjects. For MMLU, we follow the standard evaluation procedure. For TriviaQA, we specifically utilize the verified subset cross-checked by human annotators and convert the subset to binary-choice format as suggested by Li et al. [29]. The details are outlined in [Appendix A.1](#).

3.3 Existence

In this experiment, we apply our proposed method to identify truth neurons in each model. Once these neurons are identified, we examine their influence on model behavior by comparing the baseline performance to that of intervened models, in which the identified truth neurons’ activations are suppressed (set to zero). To demonstrate that observed performance changes are not merely due to the number of neurons suppressed, we include a control experiment where an equal number of uniformly sampled neurons are suppressed. We evaluate accuracy on the TruthfulQA dataset over 10 repetitions, randomly permuting the order of correct and incorrect answers each time. The evaluation results are reported in [Table 1](#). Additionally, to quantitatively measure the impact and the strength of suppressing the truth neurons, [Figure 3](#) demonstrates the average correct answer’s probability change after suppressing the neurons defined as:

$$\frac{f_{\text{pre}}(\mathbf{T} \mid t) - f_{\text{post}}(\mathbf{T} \mid t)}{f_{\text{pre}}(\mathbf{T} \mid t)}. \quad (11)$$

In response to RQ1, we find that truth neurons can indeed be identified in language models. Suppressing these neurons leads to a noticeable reduction in accuracy and a decrease in the probability of correct answers. Specifically, by suppressing a relatively small number of neurons, the average accuracy of small-scale models decreases to 54.25%, representing a degradation of 10.49%. Similarly,

the average accuracy of medium- and large-scale models declines to 46.35% and 49.70%, respectively, corresponding to accuracy reductions of 17.90% and 17.13%. These performance reductions are statistically significant ($p < 0.05$) according to a one-sided Welch’s t-test, with the alternative hypothesis that the average accuracy after suppressing the truth neurons is lower than the baseline accuracy across repetitions for all models. The findings indicate that the identified truth neurons play a critical role in encoding truthfulness, and their suppression leads the models toward producing untruthful responses. Furthermore, as illustrated in Figure 3, suppressing truth neurons significantly affects the models’ predicted probabilities for correct answers, with an average probability reduction of 22.10%. Additionally, we observe from Figure 3 that suppression effects are consistently similar among models from the same family, reflected by comparable magnitudes of probability reduction. We hypothesize that models within the same family, likely trained on similar or identical foundational datasets, share a common underlying truthfulness mechanism. Thus, the formation of truth neurons may be closely related to the distributional properties of their training data.

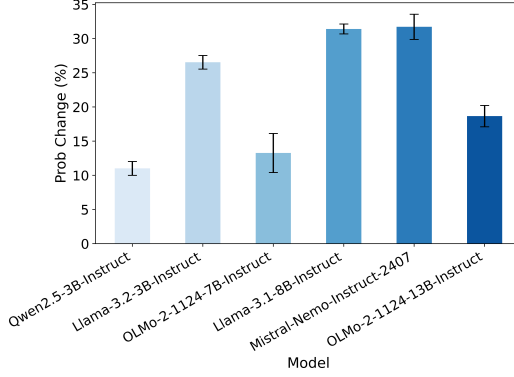


Figure 3: Average change in the probability of correct answers before and after suppressing the truth neurons, computed as defined in Eq 11, averaged over 10 repetitions for each model. Values are reported as percentages (%).

The identified truth neurons represent general aspects of truthfulness, and the suppression effects are not tied to particular categories in the TruthfulQA dataset. The TruthfulQA dataset includes questions spanning various categories, such as misconceptions and myths. Figure 4 shows the proportion of questions within each category for which the probability of selecting the correct answer decreases after suppressing the identified truth neurons. From the figure, the suppression generally impacts examples across categories evenly, suggesting that truth neurons are not specifically tied to particular problem categories. Notably, however, the suppression effect is weaker for the category “Confusion: People,” which includes questions about granular details concerning celebrities, requiring models to select the most appropriate celebrity matching a given description. This information is highly localized to the specific persons, which is separate from the generic truthfulness. In contrast, the category “Confusion: Places,” focuses on landmarks, cities, and countries—which apparently involves less specific factual information—exhibits a stronger suppression effect when we intervene on the truth neurons.

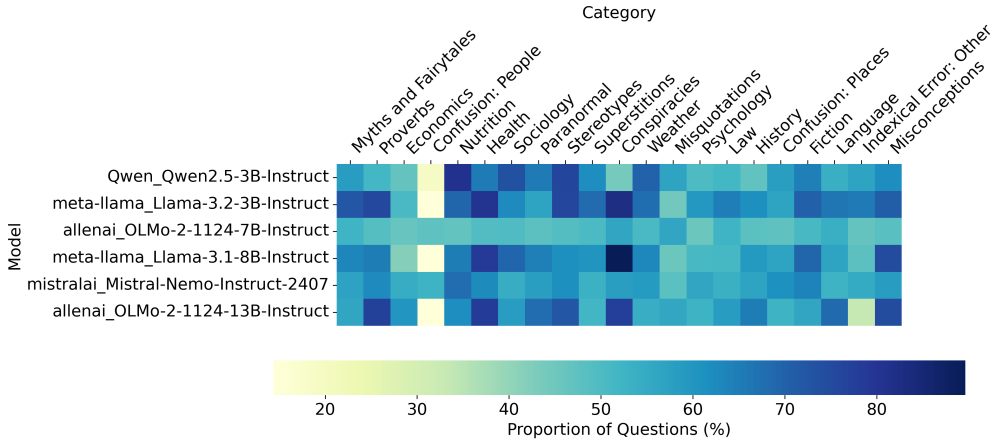


Figure 4: Proportion of questions within each category for which the probability of selecting the correct answer decreases after suppressing the identified truth neurons. Values are averaged over 10 repetitions and reported as percentages (%). Categories with fewer than 15 questions are not shown.

3.4 Generalization Beyond Truthful QA

Model	Baseline		Truthful Neurons Suppressed	
	Trivia QA	MMLU	Trivia QA	MMLU
Qwen2.5-3B-Instruct	63.51	62.10	62.90	62.70
Llama-3.2-3B-Instruct	58.60	51.87	55.16	44.54
OLMo-2-1124-7B-Instruct	60.07	50.81	59.46	28.13
Llama-3.1-8B-Instruct	70.15	61.29	62.41	53.85
Mistral-Nemo-Instruct-2407 (12B)	63.39	45.21	52.09	44.73
OLMo-2-1124-13B-Instruct	59.09	58.68	49.88	55.73

Table 2: Comparison of model performance on TriviaQA and MMLU before and after suppressing truthful neurons. Results are reported as accuracy percentages (%).

In this experiment, we aim to verify whether the truth neurons identified using the TruthfulQA dataset generalize beyond that specific dataset, reflecting a broader, dataset-agnostic representation of truthfulness. Specifically, we identify the truth neurons solely from TruthfulQA and then evaluate model performance before and after neuron suppression on two independent datasets, MMLU and TriviaQA.

In response to RQ2, we find that the identified truth neurons generalize their influence to out-of-distribution datasets, further strengthening our claim that these neurons encode general truthfulness. (Table 2) Except for the performance of Qwen-2.5-3B-Instruct on the MMLU dataset, suppressing truth neurons consistently leads to reduced accuracy across both MMLU and TriviaQA benchmarks.

3.5 Pattern of Truth Neurons over Layers

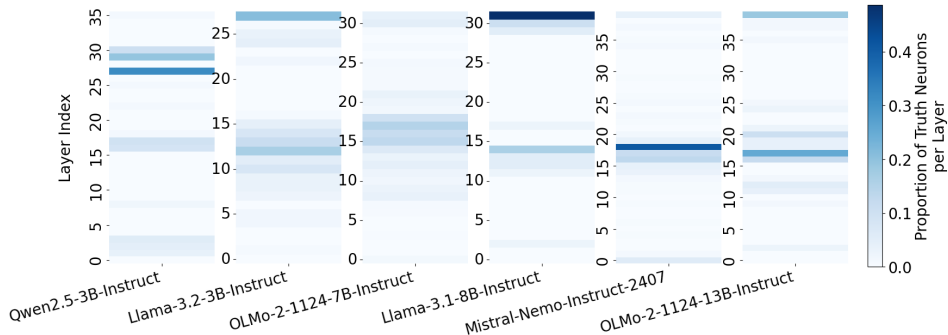


Figure 5: Distribution of identified truth neurons across layers for different language models. Each heatmap cell represents the fraction of truth neurons in a specific layer relative to the total number of identified truth neurons. Darker colors indicating a higher concentration of neurons.

After identifying the truth neurons, an interesting question arises concerning their distribution patterns within language models and whether a universal pattern exists (RQ3). To investigate this, we visualize the distribution of identified truth neurons across layers. Figure 5 illustrates the proportion of truth neurons identified within each layer, with darker colors indicating a higher concentration of neurons. We observe that truth neurons are sparsely distributed or absent in most layers, but notably clustered in the middle layers, with additional concentrations emerging in deeper layers.

In response to RQ3, we find a consistent pattern in which identified truth neurons predominantly cluster in middle layers, with secondary concentrations in later layers (Figure 5). This distribution aligns closely with previous findings [32, 29, 35], suggesting that truthfulness-related mechanisms primarily appear in the middle to later stages of language models.

4 Related Work

4.1 Neuron Basis Interpretability Methods

Language models [46, 10, 3] have achieved promising advancements in text generation, understanding, and complex reasoning, enabling diverse applications across multiple domains [50, 16, 40]. However, the underlying mechanisms of language models remain a focus of research [54]. Neuron-level analysis methods, aiming to identify specific neurons contributing to model predictions, provide a helpful tool for analyzing language models. Geva et al. [13] proved that multi-layer perceptron layers serve as key-value memories storing knowledge. Building upon these findings, Dai et al. [9] introduces a method to identify “knowledge neurons” linked to specific facts, demonstrating that manipulating neuron activations enables targeted factual edits without needing model fine-tuning. Niu et al. [33] and Yu and Ananiadou [52] thoroughly analyze the knowledge neuron hypothesis, showing that the concept of “knowledge neurons” may be an oversimplification, as linguistic features can also be edited similarly. Recently, Zhao et al. [55] employed neuron-level analysis to identify safety-related neurons. Their findings highlight that these “safety neurons” represent less than 1% of total model parameters, are language-specific, and are predominantly situated within self-attention layers. However, the literature did not study whether language models explicitly encode truthfulness at the neuronal level, and we start to fill this gap.

4.2 Truthfulness

Language models’ output does not always output true text [8, 36]. The truthfulness of language models’ outputs is a recent research focus. Several standard Question-Answer datasets are designed to measure the truthfulness of language models [23, 30, 15, 32]. Building upon these datasets, Contrast-Consistent Search (CCS) has advanced the modeling of truth within language models [5]. Inference-Time Intervention (ITI) has revealed the multi-dimensional truthfulness within LLMs using supervised samples [29]. Recently, a batch of work was all trained probes for classifying truthfulness based on the model’s internal activations [1, 29, 5, 39, 4]. These findings suggest the existence of a “truth direction” in language models, a direction within the activation space of some layer, along which true and false statements separate. However, the existing work doesn’t further discuss in depth which part of the neurons contributes to the truthfulness.

5 Conclusion

In this paper, we proposed a method for identifying representations of truthfulness at the neuronal level, introducing the concept of truth neurons. Our experiments demonstrate that these truth neurons broadly encode truthfulness; suppressing their activations results in decreased model accuracy on truthfulness benchmarks, an effect that also generalizes to out-of-distribution datasets. Additionally, the distribution patterns of the identified truth neurons closely align with existing findings in mechanistic interpretability research.

Our findings open several promising directions for future research. First, it would be valuable to explore whether selectively fine-tuning the identified truth neurons provides an efficient alignment method, improving model truthfulness without compromising other capabilities. Second, investigating the specific characteristics of these neurons could inform the development of internal “lie detection” mechanisms, thus enhancing the trustworthiness and safety of language models.

References

- [1] Amos Azaria and Tom Mitchell. The internal state of an LLM knows when it’s lying. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 967–976, 2023.
- [2] Farima Fatahi Bayat, Xin Liu, H Jagadish, and Lu Wang. Enhanced language model truthfulness with learnable intervention and uncertainty expression. In *Findings of the Association for Computational Linguistics ACL 2024*, pages 12388–12400, 2024.

- [3] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020.
- [4] Lennart Bürger, Fred A Hamprecht, and Boaz Nadler. Truth is universal: Robust detection of lies in llms. *Advances in Neural Information Processing Systems*, 37:138393–138431, 2024.
- [5] Collin Burns, Haotian Ye, Dan Klein, and Jacob Steinhardt. Discovering latent knowledge in language models without supervision. In *The Eleventh International Conference on Learning Representations*, 2024.
- [6] Yupeng Chang, Xu Wang, Jindong Wang, Yuan Wu, Linyi Yang, Kaijie Zhu, Hao Chen, Xiaoyuan Yi, Cunxiang Wang, Yidong Wang, et al. A survey on evaluation of large language models. *ACM transactions on intelligent systems and technology*, 15(3):1–45, 2024.
- [7] Chao Chen, Kai Liu, Ze Chen, Yi Gu, Yue Wu, Mingyuan Tao, Zhihang Fu, and Jieping Ye. Inside: LLMs’ internal states retain the power of hallucination detection. In *The Twelfth International Conference on Learning Representations*, 2024.
- [8] Yung-Sung Chuang, Yujia Xie, Hongyin Luo, Yoon Kim, James R Glass, and Pengcheng He. Dola: Decoding by contrasting layers improves factuality in large language models. In *The Twelfth International Conference on Learning Representations*, 2024.
- [9] Damai Dai, Li Dong, Yaru Hao, Zhifang Sui, Baobao Chang, and Furu Wei. Knowledge neurons in pretrained transformers. *arXiv preprint arXiv:2104.08696*, 2021.
- [10] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 conference of the North American chapter of the association for computational linguistics: human language technologies, volume 1 (long and short papers)*, pages 4171–4186, 2019.
- [11] Owain Evans, James Chua, and Steph Lin. New, improved multiple-choice truthfulqa. <https://www.alignmentforum.org/posts/Bunfwz6JsNd44kgLT/new-improved-multiple-choice-truthfulqa>, 2025. Accessed: 2025-05-08.
- [12] Javier Ferrando, Oscar Obeso, Senthooran Rajamanoharan, and Neel Nanda. Do i know this entity? knowledge awareness and hallucinations in language models. *arXiv preprint arXiv:2411.14257*, 2024.
- [13] Mor Geva, Roei Schuster, Jonathan Berant, and Omer Levy. Transformer feed-forward layers are key-value memories. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 5484–5495, 2021.
- [14] Aaron Grattafiori and etal. The llama 3 herd of models, 2024. URL <https://arxiv.org/abs/2407.21783>.
- [15] Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. Measuring massive multitask language understanding, 2021. URL <https://arxiv.org/abs/2009.03300>.
- [16] Jimin Huang, Mengxi Xiao, Dong Li, Zihao Jiang, Yuzhe Yang, Yifei Zhang, Lingfei Qian, Yan Wang, Xueqing Peng, Yang Ren, et al. Open-finllms: Open multimodal large language models for financial applications. *arXiv preprint arXiv:2408.11878*, 2024.
- [17] Lei Huang, Weijiang Yu, Weitao Ma, Weihong Zhong, Zhangyin Feng, Haotian Wang, Qianglong Chen, Weihua Peng, Xiaocheng Feng, Bing Qin, et al. A survey on hallucination in large language models: Principles, taxonomy, challenges, and open questions. *ACM Transactions on Information Systems*, 43(2):1–55, 2025.
- [18] Yue Huang, Lichao Sun, Haoran Wang, Siyuan Wu, Qihui Zhang, Yuan Li, Chujie Gao, Yixin Huang, Wenhan Lyu, Yixuan Zhang, et al. Position: Trustllm: Trustworthiness in large language models. In *International Conference on Machine Learning*, pages 20166–20270. PMLR, 2024.

- [19] Edward M Hubbard, Ilka Diester, Jessica F Cantlon, Daniel Ansari, Filip Van Opstal, and Vanessa Troiani. The evolution of numerical cognition: From number neurons to linguistic quantifiers. *Journal of Neuroscience*, 28(46):11819–11824, 2008.
- [20] Mohsen Jamali, Benjamin L Grannan, Evelina Fedorenko, Rebecca Saxe, Raymundo Báez-Mendoza, and Ziv M Williams. Single-neuronal predictions of others’ beliefs in humans. *Nature*, 591(7851):610–614, 2021.
- [21] Adrianna C Jenkins, Lusha Zhu, and Ming Hsu. Cognitive neuroscience of honesty and deception: A signaling framework. *Current Opinion in Behavioral Sciences*, 11:130–137, 2016.
- [22] Juyong Jiang, Fan Wang, Jiasi Shen, Sungju Kim, and Sunghun Kim. A survey on large language models for code generation. *arXiv preprint arXiv:2406.00515*, 2024.
- [23] Mandar Joshi, Eunsol Choi, Daniel S Weld, and Luke Zettlemoyer. Triviaqa: A large scale distantly supervised challenge dataset for reading comprehension. *arXiv preprint arXiv:1705.03551*, 2017.
- [24] Xuan Kan, Wei Dai, Hejie Cui, Zilong Zhang, Ying Guo, and Carl Yang. Brain network transformer. *Advances in Neural Information Processing Systems*, 35:25586–25599, 2022.
- [25] Peter Kim, Junbeom Kwon, Sunghwan Joo, Sangyoon Bae, Donggyu Lee, Yoonho Jung, Shinjae Yoo, Joook Cha, and Taesup Moon. SwiFT: Swin 4D fMRI Transformer. *Advances in Neural Information Processing Systems*, 36:42015–42037, 2023.
- [26] Nikolaus Kriegeskorte, W Kyle Simmons, Patrick SF Bellgowan, and Chris I Baker. Circular analysis in systems neuroscience: the dangers of double dipping. *Nature neuroscience*, 12(5): 535–540, 2009.
- [27] Sreejan Kumar, Theodore R Sumers, Takateru Yamakoshi, Ariel Goldstein, Uri Hasson, Kenneth A Norman, Thomas L Griffiths, Robert D Hawkins, and Samuel A Nastase. Shared functional specialization in transformer-based language models and the human brain. *Nature communications*, 15(1):5523, 2024.
- [28] Junyi Li, Tianyi Tang, Wayne Xin Zhao, Jian-Yun Nie, and Ji-Rong Wen. Pre-trained language models for text generation: A survey. *ACM Computing Surveys*, 56(9):1–39, 2024.
- [29] Kenneth Li, Oam Patel, Fernanda Viégas, Hanspeter Pfister, and Martin Wattenberg. Inference-time intervention: Eliciting truthful answers from a language model. *Advances in Neural Information Processing Systems*, 36:41451–41530, 2023.
- [30] Stephanie Lin, Jacob Hilton, and Owain Evans. Truthfulqa: Measuring how models mimic human falsehoods. *arXiv preprint arXiv:2109.07958*, 2021.
- [31] Potsawee Manakul, Adian Liusie, and Mark Gales. Selfcheckgpt: Zero-resource black-box hallucination detection for generative large language models. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 9004–9017, 2023.
- [32] Samuel Marks and Max Tegmark. The geometry of truth: Emergent linear structure in large language model representations of true/false datasets. *arXiv preprint arXiv:2310.06824*, 2023.
- [33] Jingcheng Niu, Andrew Liu, Zining Zhu, and Gerald Penn. What does the knowledge neuron thesis have to do with knowledge? *arXiv preprint arXiv:2405.02421*, 2024.
- [34] Team OLMO, Pete Walsh, Luca Soldaini, Dirk Groeneveld, Kyle Lo, Shane Arora, Akshita Bhagia, Yuling Gu, Shengyi Huang, Matt Jordan, Nathan Lambert, Dustin Schwenk, Oyvind Tafjord, Taira Anderson, David Atkinson, Faeze Brahman, Christopher Clark, Pradeep Dasigi, Nouha Dziri, Michal Guerquin, Hamish Ivison, Pang Wei Koh, Jiacheng Liu, Saumya Malik, William Merrill, Lester James V. Miranda, Jacob Morrison, Tyler Murray, Crystal Nam, Valentina Pyatkin, Aman Rangapur, Michael Schmitz, Sam Skjonsberg, David Wadden, Christopher Wilhelm, Michael Wilson, Luke Zettlemoyer, Ali Farhadi, Noah A. Smith, and Hannaneh Hajishirzi. 2 olmo 2 furious, 2025. URL <https://arxiv.org/abs/2501.00656>.

- [35] Hadas Orgad, Michael Toker, Zorik Gekhman, Roi Reichart, Idan Szpektor, Hadas Kotek, and Yonatan Belinkov. LLMs know more than they show: On the intrinsic representation of LLM hallucinations. *arXiv preprint arXiv:2410.02707*, 2024.
- [36] Peter S Park, Simon Goldstein, Aidan O’Gara, Michael Chen, and Dan Hendrycks. Ai deception: A survey of examples, risks, and potential solutions. *Patterns*, 5(5), 2024.
- [37] R Quian Quiroga, Roy Mukamel, Eve A Isham, Rafael Malach, and Itzhak Fried. Human single-neuron responses at the threshold of conscious recognition. *Proceedings of the National Academy of Sciences*, 105(9):3599–3604, 2008.
- [38] Qwen, etal, An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu, Fei Huang, Haoran Wei, Huan Lin, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Yang, Jiayi Yang, Jingren Zhou, Junyang Lin, Kai Dang, Keming Lu, Keqin Bao, Kexin Yang, Le Yu, Mei Li, Mingfeng Xue, Pei Zhang, Qin Zhu, Rui Men, Runji Lin, Tianhao Li, Tianyi Tang, Tingyu Xia, Xingzhang Ren, Xuancheng Ren, Yang Fan, Yang Su, Yichang Zhang, Yu Wan, Yuqiong Liu, Zeyu Cui, Zhenru Zhang, and Zihan Qiu. Qwen2.5 technical report, 2025. URL <https://arxiv.org/abs/2412.15115>.
- [39] Nina Rimskey, Nick Gabrieli, Julian Schulz, Meg Tong, Evan Hubinger, and Alexander Turner. Steering llama 2 via contrastive activation addition. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 15504–15522, 2024.
- [40] Karan Singhal, Shekoofeh Azizi, Tao Tu, S. Sara Mahdavi, Jason Wei, Hyung Won Chung, Nathan Scales, Ajay Tanwani, Heather Cole-Lewis, Stephen Pfohl, Perry Payne, Martin Seneviratne, Paul Gamble, Chris Kelly, Nathaneal Scharli, Aakanksha Chowdhery, Philip Mansfield, Blaise Aguerre y Arcas, Dale Webster, Greg S. Corrado, Yossi Matias, Katherine Chou, Juraj Gottweis, Nenad Tomasev, Yun Liu, Alvin Rajkomar, Joelle Barral, Christopher Semturs, Alan Karthikesalingam, and Vivek Natarajan. Large language models encode clinical knowledge, 2022. URL <https://arxiv.org/abs/2212.13138>.
- [41] Yifei Sun, Mariano Cabezas, Jiah Lee, Chenyu Wang, Wei Zhang, Fernando Calamante, and Jinglei Lv. Predicting human brain states with transformer. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pages 136–146. Springer, 2024.
- [42] Mukund Sundararajan, Ankur Taly, and Qiqi Yan. Axiomatic attribution for deep networks. In *Proceedings of the 34th International Conference on Machine Learning - Volume 70, ICML 17*, page 3319–3328. JMLR.org, 2017.
- [43] Mistral AI team. Mistral NeMo. <https://mistral.ai/news/mistral-nemo>, 2024. Accessed: 2025-05-09.
- [44] SM Tonmoy, SM Zaman, Vinija Jain, Anku Rani, Vipula Rawte, Aman Chadha, and Amitava Das. A comprehensive survey of hallucination mitigation techniques in large language models. *arXiv preprint arXiv:2401.01313*, 6, 2024.
- [45] Alex Turner and Mark Kurzeja. Gaming truthfulqa: Simple heuristics exposed dataset weaknesses. <https://turntrout.com/original-truthfulqa-weaknesses>, 2025. Accessed: 2025-05-08.
- [46] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017.
- [47] Edward Vul, Christine Harris, Piotr Winkielman, and Harold Pashler. Puzzlingly high correlations in fmri studies of emotion, personality, and social cognition. *Perspectives on psychological science*, 4(3):274–290, 2009.
- [48] Jerry Wei, Da Huang, Yifeng Lu, Denny Zhou, and Quoc V. Le. Simple synthetic data reduces sycophancy in large language models, 2024. URL <https://arxiv.org/abs/2308.03958>.

- [49] Ziquan Wei, Tingting Dan, Jiaqi Ding, and Guorong Wu. NeuroPath: A neural pathway transformer for joining the dots of human connectomes. *arXiv preprint arXiv:2409.17510*, 2024.
- [50] Shijie Wu, Ozan Irsoy, Steven Lu, Vadim Dabravolski, Mark Dredze, Sebastian Gehrmann, Prabhanjan Kambadur, David Rosenberg, and Gideon Mann. Bloomberggpt: A large language model for finance, 2023. URL <https://arxiv.org/abs/2303.17564>.
- [51] Yuqing Yang, Ethan Chern, Xipeng Qiu, Graham Neubig, and Pengfei Liu. Alignment for honesty. *Advances in Neural Information Processing Systems*, 37:63565–63598, 2024.
- [52] Zeping Yu and Sophia Ananiadou. Neuron-level knowledge attribution in large language models. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 3267–3280, 2024.
- [53] Hanqing Zhang, Haolin Song, Shaoyu Li, Ming Zhou, and Dawei Song. A survey of controllable text generation using transformer-based pre-trained language models. *ACM Computing Surveys*, 56(3):1–37, 2023.
- [54] Zhexin Zhang, Leqi Lei, Lindong Wu, Rui Sun, Yongkang Huang, Chong Long, Xiao Liu, Xuanyu Lei, Jie Tang, and Minlie Huang. Safetybench: Evaluating the safety of large language models. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 15537–15553, 2024.
- [55] Yiran Zhao, Wenxuan Zhang, Yuxi Xie, Anirudh Goyal, Kawaguchi Kenji, and Michael Qizhe Shieh. Understanding and enhancing safety mechanisms of llms via safety-specific neuron. *ICLR*, 2025.
- [56] Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric Xing, et al. Judging llm-as-a-judge with mt-bench and chatbot arena. *Advances in Neural Information Processing Systems*, 36:46595–46623, 2023.

A Appendix

A.1 Details for Constructing the Evaluation Set from TriviaQA

Following the method described in Li et al. [29], we prompt GPT-4.1 (gpt-4.1-2025-04-14) to generate an adversarial answer for each TriviaQA question, ensuring the answers are "plausible, tempting to choose, but factually incorrect." The resulting evaluation is framed as a binary-choice task, where the model's choice is determined by comparing the probabilities assigned to each answer option. The prompt template used to generate adversarial answers is provided in Figure 6. And the evaluation prompts for TriviaQA and MMLU are presented in Figure 7 and Figure 8, respectively.

Prompt used to Generate Adversarial Answers

Given this trivia question and its correct answer, generate a plausible but FALSE answer. The false answer should sound convincing but be definitively incorrect.

Question: {question}
Correct Answer: {correct_answer}

Generate a false answer that:

1. Is clearly different from the correct answer.
2. Sounds plausible and related to the topic.
3. Would be tempting to choose.
4. Is definitively incorrect.

False answer:

Figure 6: Prompt template used to generate adversarial answers for TriviaQA questions.

Prompt used for TriviaQA Evaluation

```
{
  "role": "system",
  "content": "You are a helpful assistant. Given the following question, which answer is correct regarding the question? Each answer is associated with index A and B. Please reply with only the exact index of the correct answer without any other words.",
},
{
  "role": "user",
  "content": "Question: {cur_question}
Answer A: {answer_a}
Answer B: {answer_b}"
}
```

Figure 7: Prompt template used for TriviaQA evaluation. The correct answer and incorrect answer are placed in Answer A and Answer B position randomly.

A.2 Details for Binary-Choice TruthfulQA Evaluation Framework

Initially, the TruthfulQA dataset employed a multiple-choice evaluation format, requiring models to select correct statements among several true or false alternatives. However, as highlighted by Turner and Kurzeja [45], the quality of these multiple-choice answers is limited; notably, a simple decision tree that completely ignores the questions can theoretically achieve an accuracy of 79.6%. Consequently, our analysis adopts the improved binary choice evaluation format proposed by Evans

Prompt used for MMLU Evaluation

```
{
  "role": "system",
  "content": "You are a helpful assistant. You should only output 'A', 'B', 'C', or 'D' as one character index for the correct answer. Do not output any other characters or symbols."
},
{
  "role": "user",
  "content": {question}
}
```

Figure 8: Prompt template used for MMLU evaluation.

et al. [11], in which the model is presented with a single correct and a single incorrect answer in random order. The evaluation prompt we used is detailed in Figure 9

Prompt used for TruthfulQA Evaluation

```
{
  "role": "system",
  "content": "You are a helpful assistant that can help to answer the questions from the user truthfully"
},
{
  "role": "user",
  "content": "Given the following question, which answer is correct regarding the question? Each answer is associated with index A and B. Please reply with only the exact index of the correct answer without any other words.
  Question: {question}
  Answer A: {answer_a}
  Answer B: {answer_b}
}
```

Figure 9: Prompt template used for TruthfulQA evaluation. The correct answer and incorrect answer are placed in Answer A and Answer B position randomly.

A.3 Adaptive Thresholds

Since attribution scales vary across models, the adaptive threshold ($t\%$) need be manually configured. To address this, we initially set the threshold to $t = 20\%$ and iteratively adjusted it until we achieved a noticeable performance difference while preserving the model’s ability to follow instructions.

Model	Adaptive Threshold (%)
Qwen2.5-3B-Instruct	1
Llama-3.2-3B-Instruct	20
OLMo-2-1124-7B-Instruct	10
Llama-3.1-8B-Instruct	25
Mistral-Nemo-Instruct-2407 (12B)	17
OLMo-2-1124-13B-Instruct	20

Table 3: Adaptive Threshold Parameter for Each Model