

Analyzing Key Neurons in Large Language Models

Anonymous ACL submission

Abstract

Large Language Models (LLMs) possess vast amounts of knowledge within their parameters, prompting research into methods for locating and editing this knowledge. Previous investigations have primarily focused on fill-in-the-blank tasks and locating entity-related (*usually single-token facts*) information in relatively small-scale language models. However, several key questions remain unanswered: (1) *How can we effectively locate query-relevant neurons in contemporary autoregressive LLMs, such as LLaMA and Mistral?* (2) *How can we address the challenge of long-form text generation?* (3) *Are there localized knowledge regions in LLMs?* In this study, we introduce Neuron Attribution-Inverse Cluster Attribution (NA-ICA), a novel architecture-agnostic framework capable of identifying key neurons in LLMs. NA-ICA allows for the examination of long-form answers beyond single tokens by employing the proxy task of multi-choice question answering. To evaluate the effectiveness of our detected key neurons, we construct two multi-choice QA datasets spanning diverse domains and languages. Empirical evaluations demonstrate that NA-ICA outperforms baseline methods significantly. Moreover, analysis of neuron distributions reveals the presence of visible localized regions, particularly within different domains. Finally, we demonstrate the potential applications of our detected key neurons in knowledge editing and neuron-based prediction.

1 Introduction

Large Language Models (LLMs) contain substantial amounts of knowledge within their parameters. Existing research endeavors to locate and edit this knowledge through gradient-based methods (Dai et al., 2022) or causality-based methods (Meng et al., 2022a). These methods typically employ fill-in-the-blank tasks, such as “Paris is the capital of ___”, to ascertain the correlation

between the query and neurons or layers in the Feed-forward Networks (FFNs) of BERT (Kenton and Toutanova, 2019) and GPT (Radford et al.). Another branch of pioneering research attempts to locate functional regions in small-size language models such as BERT and GPT-small, including linguistic regions (Zhang et al., 2024b), factual subnetworks (Ren and Zhu, 2022; Bayazit et al., 2023), and modular structures (Zhang et al., 2023; Conmy et al., 2023).

While these studies successfully analyze the internal reasoning behaviors of LLMs, three significant questions remain underexplored: (1) How can we effectively locate query-relevant neurons in contemporary autoregressive LLMs, such as LLaMA (Touvron et al., 2023) and Mistral (Jiang et al., 2023), given that their FFNs architectures differ from those of BERT and GPT? (2) How can we address the challenge of long-form text generation, as previous methods have been limited to single-token entity facts? (3) Are there localized knowledge regions in LLMs analogous to the localized functional regions observed in human brains (Brett et al., 2002)?

To address the first two questions, we introduce a novel framework named *Neuron Attribution-Inverse Cluster Attribution (NA-ICA)* designed to identify key neurons in LLMs. The principal advantages of NA-ICA are its architecture-agnostic nature and its capability of handling long-form text generation effectively. The overall structure of the framework is depicted in Figure 1. NA-ICA draws inspiration from the TF-IDF keyword extraction method (Salton, 1983), aiming to extract significant neurons for each input query. The process begins by transforming an open-ended generation task into a multiple-choice question-answering format. By employing prompt engineering, we constrain LLMs to generate only the option letter rather than the complete answer. This approach allows for the examination of long-form generation beyond sin-

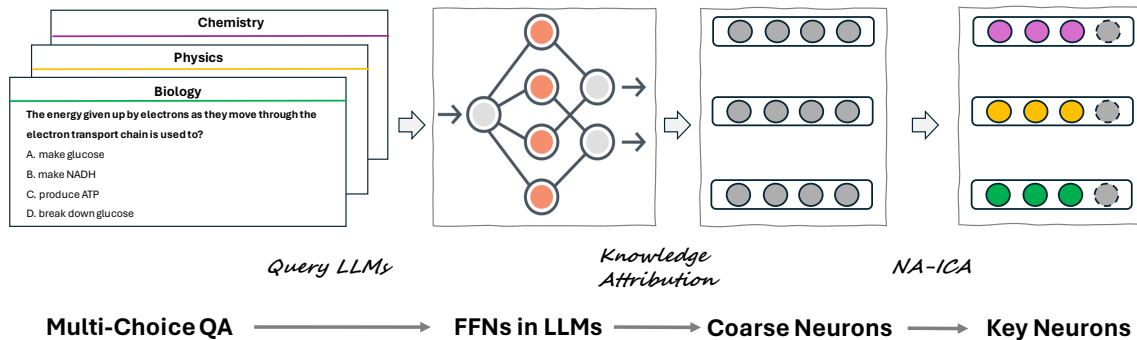


Figure 1: The overall framework of our Attribution-Inverse Cluster Attribution (NA-ICA) which aims to detect query-related key neurons. Neurons with solid lines mean key neurons while dashed ones mean common neurons that are shared across different queries.

gle tokens and extends previous methodologies to autoregressive LLMs. Subsequently, we adapt the Knowledge Attribution method (Dai et al., 2022) to compute *Neuron Attribution*, which elucidates the relationship between neurons and the input query. We then gather clusters for a series of queries and calculate the *Inverse Cluster Attribution*. This step mitigates the influence of neurons that recur across clusters (or queries). The final step involves multiplying the neuron attribution and inverse cluster attribution values to pinpoint key neurons. Additionally, we identify certain *Common Neurons* that are associated with common words, punctuation marks, and option letters. Excluding these common neurons enhances the detection of key neurons. Empirical evaluations demonstrate that our proposed method outperforms baseline approaches.

To investigate the existence of localised knowledge regions, we construct two multi-choice QA datasets encompassing various domains and languages. Then, we visualize the geographical locations of the detected key neurons in LLaMA. Our findings indicate that distinct localized regions emerge in the middle layers, particularly for domain-specific neurons. Language neurons are more sparse but show a certain degree of regionality. Additionally, we observed that common neurons are concentrated in the top layer, predominantly expressing frequently used tokens.

In summary, our main contributions are four-fold: (1) **A scalable method:** we propose NA-ICA to detect key neurons in LLMs; NA-ICA method is architecture-agnostic and can deal with long-form generations. (2) **Two new datasets:** we curate two multi-choice QA datasets that contain different types of knowledge, namely Domain Knowledge

and Language knowledge. (3) **In-depth studies:** we are the first to show that there are visible localized regions in LLaMA. (4) **Potential applications:** we show that NA-ICA might be useful for knowledge editing and neuron-based prediction.

2 Related Work

2.1 Locating Knowledge in LLMs

LLMs contain extensive knowledge within their parameters, encompassing factual (Petroni et al., 2019; Zhou et al., 2020; Jiang et al., 2020; Roberts et al., 2020; Pezeshkpour, 2023), linguistic (Liu et al., 2019; Jawahar et al., 2019; Chen et al., 2023) and domain-specific knowledge (Sung et al., 2021; Frieder et al., 2024). Despite this, the mechanisms and locations of knowledge storage within these models remain unclear. Recent mechanistic studies suggest that knowledge is primarily stored in the FFNs (Feed-forward Networks) layers of Transformers (Geva et al., 2021, 2022). Ongoing research is focused on developing methods to precisely identify and locate this knowledge within the FFNs layers. Given an input, *gradient-based methods* (Ancona et al., 2019; Dai et al., 2022) quantify the sensitivity of model outputs to internal model components, identifying relevant neurons. However, these studies focus exclusively on traditional neural architectures and encoder-only models like BERT, leaving decoder-only models such as GPT and LLaMA underexplored. *Causality-based methods* employ causal mediation analysis to discern the particular layers associated with a given factual input (Meng et al., 2022a). Subsequent research adopts the locate-and-edit paradigm to refine the knowledge within LLMs (Meng et al., 2022b; Ju and Zhang, 2023; Zhang et al., 2024a).

While previous approaches have effectively identified specific information in LLMs, they commonly rely on the fill-in-the-blank cloze task to evaluate the factual capabilities of language models. For instance, they use a prompt query like “Paris is the capital of ___” to locate weights associated with the France entity. However, this methodology has limited applicability, as language models exhibit the capacity to generate long-form and open-ended responses to diverse queries. In contrast to prior methodologies, our approach leverages the proxy task of multiple-choice QA for knowledge localization. This alternative strategy renders the localization process architecture-agnostic and facilitates the handling of long-form content generation.

2.2 Analyzing Knowledge Distribution in LLMs

Given the human-like reasoning capabilities observed in LLMs across various tasks (Zhao et al., 2023), and since our brain contains functional locations associated with distinct cognitive processes (Brett et al., 2002; Bjaalie, 2002; Gholipour et al., 2007), we ask whether there are similar regions in LLMs. Previous investigations have explored the behaviors of individual neurons indicating that a neuron can encode multiple concepts (Bolukbasi et al., 2021) while a concept can also be distributed across multiple neurons (Dalvi et al., 2019; Durrani et al., 2020; Chen et al., 2024). Subsequent endeavors have sought to identify functional regions in LLMs, encompassing linguistic regions (Zhang et al., 2024b), factual subnetworks (Ren and Zhu, 2022; Bayazit et al., 2023), and modular structures (Zhang et al., 2023; Conmy et al., 2023). These studies have systematically investigated localized behaviors in smaller-scale language models, such as BERT and GPT-small. Building upon these foundations, our research embarks on the examination of knowledge locations in larger-size LLMs, specifically those with 7B parameters, spanning multiple knowledge domains.

3 Background

Feed-forward Networks in LLMs Feed-forward networks (FFNs) are widely used by transformer-based language models. Geva et al. (2021) reveal that FFNs emulate key-value memories and their outputs are responsible for refining the final output distribution over the

vocabulary. Although traditional two-layer FFNs in BERT (Kenton and Toutanova, 2019) and GPT-2 (Radford et al.) have been studied well, the behaviors of FFNs in modern LLMs such as LLaMA (Touvron et al., 2023), Mistral (Jiang et al., 2023), and Gemma (Team et al., 2024) are not well-explored. These LLMs adopt Gated Linear Units (GLUs) (Dauphin et al., 2017) in their FFNs, which can be formulated as follows:

$$\text{FFN}(\mathbf{X}) = (\mathbf{X}\mathbf{W}^U \odot \text{SiLU}(\mathbf{X}\mathbf{W}^G)) \mathbf{W}^D \quad (1)$$

Here, $\mathbf{X} \in \mathbb{R}^{n \times d}$ is the input sequence, n is the number of tokens and d is the dimension of input vectors; $\mathbf{W}^U \in \mathbb{R}^{d \times m}$, $\mathbf{W}^G \in \mathbb{R}^{d \times m}$, $\mathbf{W}^D \in \mathbb{R}^{m \times d}$ are parameter matrices and \odot is the Hadamard product; finally SiLU (Elfwing et al., 2018) is the activation function.

Knowledge Neurons Dai et al. (2022) propose a gradient-based *Knowledge Attribution* to identify the knowledge neurons in BERT by using the fill-in-the-blank cloze task. Their method evaluates the contribution of each neuron in FFNs to the knowledge predictions. Given a prompt q “Paris is the capital of ___”, the probability of the correct answer predicted by a language model can be formulated as:

$$P_q(\hat{w}_i^l) = p(y^* | x, w_i^l = \hat{w}_i^l) \quad (2)$$

where y^* is the correct answer (France); w_i^l denotes the i -th intermediate neuron in the l -th layer in FFNs; \hat{w}_i^l is a constant we assign to w_i^l .

In order to measure the attribution score (or contribution) of a neuron, they gradually change the w_i^l from 0 to its original value computed during the forward pass through the LLM and integrate the gradients (Sundararajan et al., 2017):

$$\text{Attr}(w_i^l) = \bar{w}_i^l \int_{\alpha=0}^1 \frac{\partial P_q(\alpha \bar{w}_i^l)}{\partial w_i^l} d\alpha \quad (3)$$

where $\frac{\partial P_q(\alpha \bar{w}_i^l)}{\partial w_i^l}$ is the gradient with regard to w_i^l . $\text{Attr}(\cdot)$ accumulates the output probability change as α gradually varies from 0 to 1. The attribution measures the contribution of the neuron w_i^l to the correct answer. In practice, the score is estimated by using Riemann Approximation:

$$\hat{\text{Attr}}(w_i^l) = \frac{\bar{w}_i^l}{m} \sum_{k=1}^m \frac{\partial P_q(\frac{k}{m} \bar{w}_i^l)}{\partial w_i^l} \quad (4)$$

where m is the number of the estimation steps. Finally, they identify a coarse set of knowledge neurons whose attribution scores are greater than a threshold t .

4 Locating Key Neurons in Autoregressive LLMs

While Knowledge Attribution (Dai et al., 2022) effectively identifies neurons linked to factual queries, its applicability is limited to Encoder-only architectures, and it mandates the output to be a single-token word. In response to these constraints, we propose a simple yet effective pipeline named **Attribution-Inverse Cluster Attribution (NA-ICA)**, which is architecture-agnostic and capable of handling long-form generation. The overall framework is shown in Figure 1. NA-IQA first resorts to the proxy task of multi-choice QA to deal with long-form answers. Subsequently, the framework extracts key neurons for each query using our designed NA-ICA score.

4.1 Multi-Choice QA Transformation

Given the biological question “*The energy given up by electrons as they move through the electron transport chain is used to?*”, the correct answer can be the long-form text “*produce ATP*”. To deal with long-form answers, we advocate for the transformation of questions and their corresponding answers into a multiple-choice framework, as illustrated in Figure 1. This approach involves the generation of incorrect options by randomly sampling answers within the same domain. Following this, the LLM is prompted to produce only the option letter. Subsequently, we investigate the key knowledge neurons correlated with the input query. To mitigate the impact of randomness, we devise multiple prompt templates and systematically shuffle the order of options to prevent the model from learning spurious correlations based on option letters. These prompt templates are detailed in Table A1.

4.2 Neuron Attribution-Inverse Cluster Attribution

In our pursuit of locating neurons associated with specific queries, we compute the score of NA-ICA for each neuron, drawing inspiration from the principles of TF-IDF (Salton, 1983) for keyword extraction. Beginning with a given query, NA-ICA employs *neuron attribution* to derive a coarse group key neurons, termed as *clusters*. Each

neuron within this cluster is assigned an attribution score indicative of its relevance to the query, akin to the computation of term frequency. Given our objective of identifying critical neurons closely correlated with their respective queries, we use *inverse cluster attribution* to filter out neurons shared across different clusters (or queries). Finally, we find some neurons appear across multiple clusters, embodying common knowledge or sense, which we denote as *Common Neurons*. Further refinement of key neuron extraction involves the exclusion of these common neurons, which can enhance the precision of identifying critical neural correlates.

Neuron Attribution To extend our methodology to Gated Linear Units (GLUs), which comprise two linear transformations followed by a gating mechanism, we adapt the Knowledge Attribution approach (Eq 5). In GLUs, the linear transformations involve computing a linear combination of input features, denoted by $f = \mathbf{XW}^U$. Additionally, the gating mechanism, represented by $g = \text{SiLU}(\mathbf{XW}^G)$, determines the extent to which each input component should be forwarded, thereby enabling the model to emphasize important features while suppressing irrelevant ones. To compute the relevant attribution, we can use either $\frac{\partial P_q}{\partial f}$ or $\frac{\partial P_q}{\partial g}$ and we choose to use the former since our empirical study shows it can obtain better key neurons (see details in the Table A3). Given a query q , instantiation using our templates yields a query set $\mathcal{Q} = \{q_1, q_2, \dots, q_{|\mathcal{Q}|}\}$, and the attribution score of the neuron n_i^l can be denoted as:

$$\text{na}(n_i^l) = \sum_{j=1}^{|\mathcal{Q}|} \frac{f_i^l}{m} \sum_{k=1}^m \frac{\partial P_{q_j}(\frac{k}{m} f_i^l)}{\partial f_i^l} \quad (5)$$

Here, we sum up the scores of different instantiated templates together as the final attribution score.

Inverse Cluster Attribution With the attribution score, we can obtain a list of coarse clusters for each query $\mathcal{C} = \{c_1, c_2, \dots, c_{|\mathcal{C}|}\}$, where c is a cluster that consists of neurons whose attribution score is higher than some threshold t . The frequent appearance of some neurons across queries of different fields reveals that they are not critical neurons to the input query. To decrease their impact, we calculate the inverse cluster attribution:

$$\text{ica}(n_i^l) = \log \frac{|\mathcal{C}|}{|\{c : c \in \mathcal{C} \text{ and } n_i^l \in c\}| + 1} \quad (6)$$

Common Neurons We observe that some neurons with a relatively high NA-ICA score are still shared across clusters. Through case studies (as shown in Table 4), we demonstrate that they express commonly used concepts such as option letters (“A” and “B”) or stop words (“and” and “the”). Therefore, we count the frequency of each neuron across clusters. If the frequency is higher than the $u\%$ of total clusters, we assign the given neuron into the common neuron set.

Key Neurons Given a query, the NA-ICA of a neuron can be computed as :

$$\text{naica}(n_i^l) = \text{na}(n_i^l) \times \text{ica}(n_i^l) \quad (7)$$

We select top- v neurons with the highest score from the detected cluster and further remove common neurons to refine the key neuron set.

5 Analyzing Detected Key Neurons

5.1 Dataset Construction

We construct two datasets to locate knowledge neurons that cover two different categories: *subject domains and languages*.

Domain Dataset is derived from MMLU (Hendrycks et al., 2020), a multiple-choice QA benchmark designed to evaluate models across a wide array of subjects with varying difficulty levels. The subjects encompass traditional disciplines such as mathematics and history, as well as specialized fields like law and ethics. In our study, we select six high school exam subjects from the test set: *Biology, Physics, Chemistry, Mathematics, Computer Science, and Geography*.

Language Dataset is adapted from Multilingual LAMA (Kassner et al., 2021), which is a dataset to investigate knowledge in language models in a multilingual setting covering 53 languages. We select six languages for the birth_place relation: *Arabic, English, French, Japanese, Russian and Chinese*.

To mitigate sensitivity to prompts and option orders, each query is instantiated with multiple distinct templates (as shown in Table A1), and the option orders are shuffled each time. The statistics of our datasets are shown in Table 1 and examples can be found in Table A2.

Domain	Bio	Phy	Chem	Math	CS	Geo	Total
Num	100	100	100	100	52	100	552
Language	Ar	En	Fr	Ja	Ru	Zh	Total
Num	100	100	100	100	100	100	600

Table 1: Statistics of our constructed datasets.

5.2 Baselines

We compare our NA-ICA to three other neuron-level baselines¹: **Random Neurons** are randomly selected from FFNs and we make sure they have the same number of neurons of NA-ICA; **Knowledge Neurons*** is adapted from knowledge attribution (Dai et al., 2022) by using multi-choice QA task; **NA-ICA w/ Common Neurons** is a variant without removing common neurons.

5.3 Experimental Settings

We mainly study the knowledge neurons in LLaMA-7B (Touvron et al., 2023) and we use the instruction-tuned version so that the model is more responsive to our prompts. LLaMA-7B consists of 32 layers with the FFN hidden dimension of 11008. Besides, we also conduct experiments for Mistral-7B (Jiang et al., 2023) to validate whether that our method can obtain consistent findings over different models. Note that our framework can be easily extended to larger-size LLMs.

As for the used hyper-parameters, the number of estimation steps was set to $m = 16$ and the attribution threshold t to 0.2 times the maximum attribution score. The template number was $|\mathcal{Q}| = 3$, the frequency u for obtaining common neurons was 30%, and the top- v for select key neurons was 20. We ran all experiments on three NVIDIA-V100. It took 120 seconds on average to locate neurons for a query with three prompt templates.

5.4 Statistics of Detected Key Neurons

Table 3 presents the number of detected key neurons for each domain and language, averaging between 12 and 17 neurons. Figure 2a illustrates the overlap rates among different domains and languages. It is evident that domains exhibit higher overlap rates compared to languages, reflecting interconnected and interdisciplinary nature. For instance, the overlap rate between biology and geography is 0.49, attributable to fields like biogeography, which examines the distribution of species

¹We do not compare to ROME (Meng et al., 2022a) since it locates layers instead of neurons

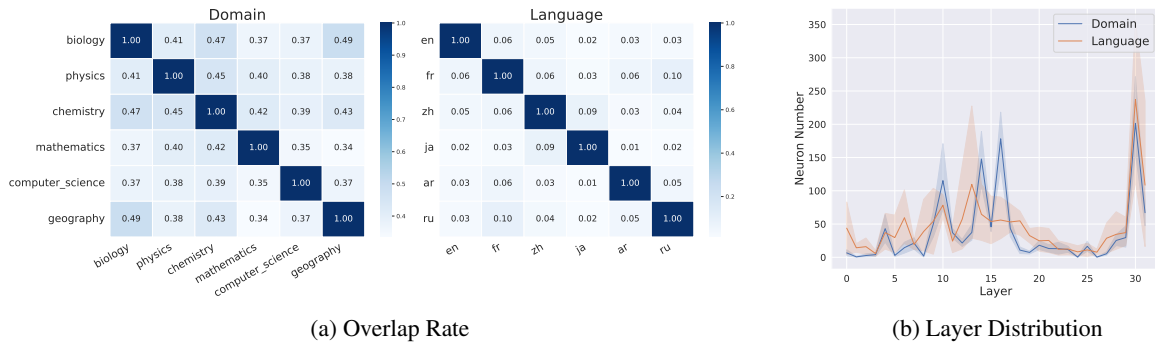


Figure 2: Overlap rates and distributions of found key neurons.

Model	Domain						Language					
	Boost			Suppress			Boost			Suppress		
	↑ Related	↑ Unrelated	Ratio	↓ Related	↓ Unrelated	Ratio	↑ Related	↑ Unrelated	Ratio	↓ Related	↓ Unrelated	Ratio
Random Neurons	-0.03	-0.03	1.0	+0.06	+0.11	0.55	+0.08	+0.04	2.0	-0.01	-0.01	1.0
Knowledge Neurons* (2022)	+932.05	+921.84	1.0	-85.70	-85.34	1.0	+1081.33	+161.98	6.7	-86.74	-48.18	1.8
NA-ICA w/ Common Neurons	+919.03	+328.49	2.8	-59.34	-33.59	1.8	+606.54	+54.84	10.4	-71.45	-8.40	8.5
NA-ICA	+77.23	+17.55	4.4	-27.65	-4.95	5.6	+218.03	+5.20	41.9	-54.64	+3.71	15.2

Table 2: Average probability percentage changes of the correct answers by boosting (↑) or suppressing (↓) the key neurons. The Ratio metric is calculated by $\frac{|\text{Related}|}{|\text{Unrelated}|}$, and a bigger value shows a higher impact of the detected neurons. The LLM here is LLaMA-7B (Touvron et al., 2023)

Domain	Bio	Phy	Chem	Math	CS	Geo	Avg
Num	13.1	13.3	12.8	11.1	14.3	12.7	12.9
Language	Ar	En	Fr	Ja	Ru	Zh	Avg
Num	12.4	14.4	12.7	16.6	15.8	15.0	14.5

Table 3: Average number of key neurons.

and ecosystems in geographic space. Regarding layer distribution, the key neurons are predominantly located in the middle layers (15-18) and the top layers (around 30), as depicted in Figure 2b.

5.5 Key Neurons Can Impact the Prediction

To validate the impact of our identified key neurons, we replicate the experiments by Dai et al. (2022), updating the values of key neurons using two methods: given a query and the value of f_i^l , we either (1) boost the key neurons by doubling the value $f_i^l = 2 \times f_i^l$; or (2) suppress the key neuron by making $f_i^l = 0$. For each query, we record the percentage change in the probability of the correct answer, thereby assessing the extent to which the key neurons influence the predictions of LLMs. We compare our NA-ICA approach to other baseline methods and include a control group to determine whether the same key neurons affect the predictions of randomly selected queries from unrelated fields (*Unrelated*).

Table 2 presents the overall performance of

various methods. Our NA-ICA method consistently outperforms other baselines, evidenced by its higher impact ratio. This indicates that our identified key neurons significantly affect the probability of correct answers while exerting a relatively low impact on unrelated queries. For instance, our method achieves a boosting ratio of 41.9 on the language dataset, the highest among the baselines. Additionally, common neurons affect both related and unrelated queries, and their removal results in clear performance improvements.

Furthermore, Figure 3 illustrates the percentage change in probability for each domain and language. Again, we can clearly observe the effectiveness of our detected key neurons. Additionally, we performed supplementary experiments on Mistral-7B. The results, presented in Figure A2, consistently support our conclusions.

5.6 Are There Localized Regions in LLMs?

Given our ability to identify key neurons for each query, it is intriguing to explore whether LLMs exhibit localized regions for each domain or language, analogous to the functional localizations in the human brain (Brett et al., 2002). To investigate this, we visualize domain- or language-specific neurons on a 2D geographical heatmap. The width of the heatmap corresponds to the dimension of FFNs in LLaMA-7B (11008), and the length represents

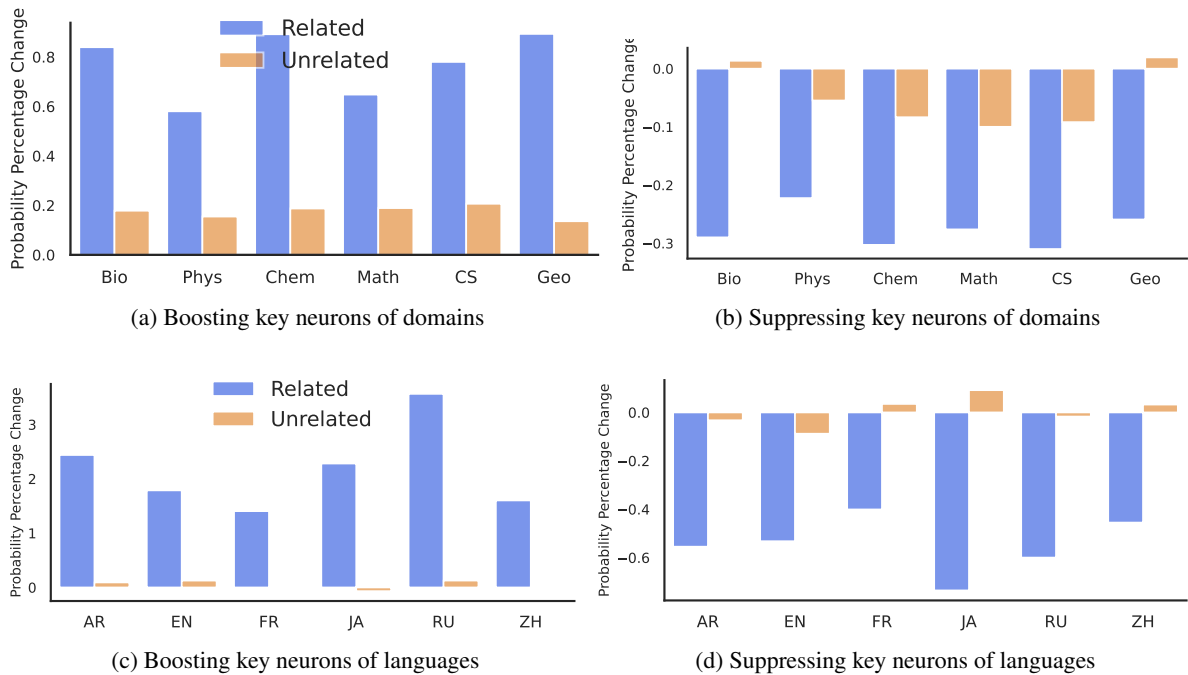


Figure 3: The correct probability percentage change across different domains and languages. The LLM here is LLaMA-7B (Touvron et al., 2023)

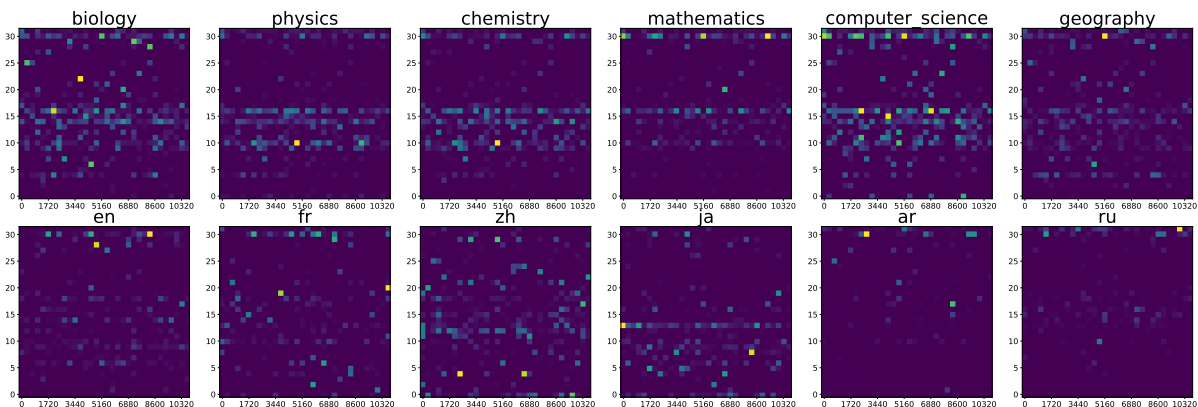


Figure 4: Geographical heatmap of detected key neurons for different domains and languages. The value is calculated by our $naica(n_i^l)$. The LLM here is LLaMA-7B (11008 × 32) (Touvron et al., 2023)

472 the layer depth (32). We accumulate the value of
 473 $naica(n_i^l)$ to populate the heatmap. Figure 4 dis-
 474 plays the geographical locations of key neurons
 475 in LLaMA-7B across various academic domains
 476 and languages. The distribution of key neurons ap-
 477 pears sparse but with distinct regions, particularly
 478 for different domains. Notably, certain regions are
 479 visible in the middle layers (10-15), suggesting spe-
 480 cific neuron patterns. In contrast, language neurons
 481 are more sparsely distributed with smaller regions,
 482 and languages like Arabic and Russian exhibit less
 483 localized properties. Apart from visualizing the ge-
 484 ographical location of key neurons, we also analyze
 485 the semantic location using their associated vector
 486 values in \mathbf{W}^D . Our findings suggest that there are
 487 no apparent clusters across different domains, as

488 these values likely represent intermediate states in a
 489 subspace distinct from the one used for final token
 490 prediction. Future work should consider new ways
 491 to map neurons to more discriminative semantic
 492 spaces. The details are provided in Appendix A.

5.7 The Function of Common Neurons 493

494 To gain insights into the function of common
 495 neurons, we also visualize their locations within
 496 LLaMA-7B. Figure 5 shows the common neurons
 497 for the domain and language dataset. We can ob-
 498 serve that they tend to appear at the top layer. To
 499 further understand their meanings, we project the
 500 matrix \mathbf{W}^D in Equation 1 to the vocabulary space
 501 and select the top-k tokens with the highest prob-
 502 ability. Table 4 lists the predicted tokens, which

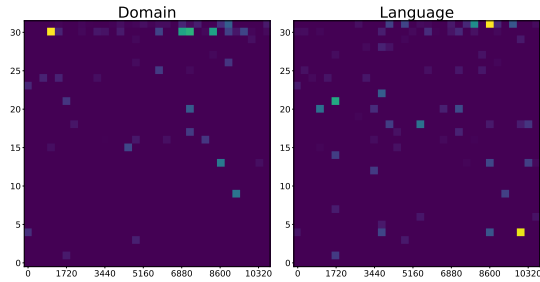


Figure 5: The distribution of common neurons.

Neuron	Top-k tokens
n_{2725}^{31}	_in, _and, _to, _for, _today, _at, _as
n_{10676}^{31}	_July, _June, _March, _April, _November
n_{10075}^{30}	., -, (, :,), [, -
n_{5202}^{31}	_respectively, _while, _and, _initially
n_{5778}^{31}	_C, C, _c, c, '_ced'
n_{7670}^{31}	_B, B, _Bill, _Bh, '_Bureau'

Table 4: Tokens predicted by the common neurons.

include common words, punctuation marks, and option letters. These findings reinforce the notion that common neurons are not critical for specific queries.

6 Potential Applications

We provide two usage examples to showcase the potential applications of our detected key neurons: *Knowledge Editing* and *Neuron-Based Prediction*.

6.1 Knowledge Editing

We adjust the values of key neurons by either boosting or suppressing them to determine if we can change the prediction of a query from incorrect to correct or vice versa. Table 5 presents the success rates of knowledge editing on our constructed language datasets. Our observations indicate that NA-ICA achieves higher success rates on related queries and lower rates on unrelated queries, demonstrating that our method outperforms other baselines. In contrast, the baseline of knowledge neurons cannot significantly differentiate related and unrelated queries.

6.2 Neuron-Based Prediction

In our second case study, we test whether the correct answers to domain-specific questions can be predicted solely based on the activity of the associated domain-specific neurons. To this end, we make predictions on multiple-choice questions by selecting the option with the overall highest gradient to the key neurons for the given domain. We experiment on a specifically constructed

Model	Boost		Suppress	
	↑ Related	↑ Unrelated	↓ Related	↓ Unrelated
Random Neurons	0.37	0.10	0.54	0.27
Knowledge Neurons (2022)	14.73	11.76	16.19	14.78
Ours	10.06	1.78	20.14	1.60

Table 5: Successful rates of knowledge editing.

Model	Biology (Acc.)	Chemistry (Acc.)	Geography (Acc.)
Random guess	0.25	0.25	0.25
Prompt-based model pred.	0.96	0.71	0.89
Neuron-based pred.	0.96	0.67	0.89

Table 6: Accuracy of neuron-based prediction on selected domains in comparison with the standard prompt-based model prediction. The LLM here is LLaMA-7B.

MMLU (Hendrycks et al., 2020) validation set with a different set of questions than those used to determine the key domain neurons (see Appendix B for details on our experimental strategy). The results are summarised in Table 6. We observe that the accuracy of the neuron-based predictions is very close to the accuracy of the prompt-based method of using the entire model (the used templates are shown in Table A1). This suggests that the activity of identified neurons can be indicative of the model’s performance on a given task. Investigating how this finding could be leveraged in applications like fact-checking and hallucination detection presents a promising line of future work.

7 Conclusion

In this study, we introduce a novel framework, NA-ICA, for identifying key neurons in contemporary autoregressive language models, such as LLaMA and Mistral. NA-ICA leverages a multi-choice QA proxy task to address the complexity of long-form answers, extending beyond simple factual entities. Meanwhile, it adopts strategies of inverse cluster attribution and common neuron removal to refine key neurons. To validate our approach, we curated two datasets encompassing diverse domains and languages. Our experimental results show that NA-ICA outperforms existing baselines in identifying query-relevant neurons. Additionally, this study pioneers the exploration of localized knowledge regions in LLMs and demonstrates the potential usages of identified key neurons in applications such as knowledge editing and neuron-based prediction. We hope that our findings are beneficial for further research in understanding the knowledge mechanisms underlying LLMs.

568 Limitations

569 In our study, we employ a multi-choice QA proxy
570 task to investigate the long-form knowledge stored
571 in LLMs. Although our framework can effectively
572 detect key neurons, future research needs to address
573 the challenge of authentic open-ended generation,
574 which remains a significant area for development.
575 Additionally, despite our efforts to eliminate com-
576 mon neurons, some neurons within the identified
577 key neuron set still correspond to option letters.
578 This indicates that our current method requires fur-
579 ther refinement to remove these spurious key neu-
580 rons. Moreover, the language dataset used in our
581 study is limited to the Birth_place relation. To
582 gain a more comprehensive understanding of mul-
583 tilingual knowledge in LLMs, future work should
584 include a broader range of relations. This expan-
585 sion will enable a more thorough investigation into
586 the diverse types of knowledge encoded in these
587 models across different languages.

588 References

589 Marco Ancona, Enea Ceolini, Cengiz Öztireli, and
590 Markus Gross. 2019. Gradient-based attribution
591 methods. *Explainable AI: Interpreting, explaining
592 and visualizing deep learning*, pages 169–191.

593 Deniz Bayazit, Negar Foroutan, Zeming Chen, Gail
594 Weiss, and Antoine Bosselut. 2023. Discovering
595 knowledge-critical subnetworks in pretrained lan-
596 guage models. *arXiv preprint arXiv:2310.03084*.

597 Jan G Bjaalie. 2002. Localization in the brain: new
598 solutions emerging. *Nature reviews neuroscience*,
599 3(4):322–325.

600 Tolga Bolukbasi, Adam Pearce, Ann Yuan, Andy Co-
601 enen, Emily Reif, Fernanda Viégas, and Martin Wat-
602 tenberg. 2021. An interpretability illusion for bert.
603 *arXiv preprint arXiv:2104.07143*.

604 Matthew Brett, Ingrid S Johnsrude, and Adrian M Owen.
605 2002. The problem of functional localization in the
606 human brain. *Nature reviews neuroscience*, 3(3):243–
607 249.

608 Lihu Chen, Gael Varoquaux, and Fabian Suchanek.
609 2023. The locality and symmetry of positional en-
610 codings. In *Findings of the Association for Com-
611 putational Linguistics: EMNLP 2023*, pages 14313–
612 14331.

613 Yuheng Chen, Pengfei Cao, Yubo Chen, Kang Liu, and
614 Jun Zhao. 2024. Journey to the center of the knowl-
615 edge neurons: Discoveries of language-independent
616 knowledge neurons and degenerate knowledge neu-
617 rons. In *Proceedings of the AAAI Conference on Ar-
618 tificial Intelligence*, volume 38, pages 17817–17825.

Arthur Conmy, Augustine Mavor-Parker, Aengus Lynch, 619
Stefan Heimersheim, and Adrià Garriga-Alonso. 620
2023. Towards automated circuit discovery for mech- 621
anistic interpretability. *Advances in Neural Informa- 622
tion Processing Systems*, 36:16318–16352. 623

Damai Dai, Li Dong, Yaru Hao, Zhifang Sui, Baobao 624
Chang, and Furu Wei. 2022. Knowledge neurons in 625
pretrained transformers. In *Proceedings of the 60th 626
Annual Meeting of the Association for Computational 627
Linguistics (Volume 1: Long Papers)*, pages 8493– 628
8502. 629

Fahim Dalvi, Nadir Durrani, Hassan Sajjad, Yonatan Be- 630
linkov, Anthony Bau, and James Glass. 2019. What 631
is one grain of sand in the desert? analyzing indi- 632
vidual neurons in deep nlp models. In *Proceedings 633
of the AAAI Conference on Artificial Intelligence*, 634
volume 33, pages 6309–6317. 635

Yann N Dauphin, Angela Fan, Michael Auli, and David 636
Grangier. 2017. Language modeling with gated con- 637
volutional networks. In *International conference on 638
machine learning*, pages 933–941. PMLR. 639

Nadir Durrani, Hassan Sajjad, Fahim Dalvi, and 640
Yonatan Belinkov. 2020. Analyzing individual neu- 641
rons in pre-trained language models. In *Proceedings 642
of the 2020 Conference on Empirical Methods in 643
Natural Language Processing (EMNLP)*, pages 4865– 644
4880. 645

Stefan Elfving, Eiji Uchibe, and Kenji Doya. 2018. 646
Sigmoid-weighted linear units for neural network 647
function approximation in reinforcement learning. 648
Neural networks, 107:3–11. 649

Simon Frieder, Luca Pinchetti, Ryan-Rhys Griffiths, 650
Tommaso Salvatori, Thomas Lukasiewicz, Philipp 651
Petersen, and Julius Berner. 2024. Mathematical ca- 652
pabilities of chatgpt. *Advances in Neural Information 653
Processing Systems*, 36. 654

Mor Geva, Avi Caciularu, Kevin Wang, and Yoav Gold- 655
berg. 2022. Transformer feed-forward layers build 656
predictions by promoting concepts in the vocabulary 657
space. In *Proceedings of the 2022 Conference on 658
Empirical Methods in Natural Language Processing*, 659
pages 30–45. 660

Mor Geva, Roei Schuster, Jonathan Berant, and Omer 661
Levy. 2021. Transformer feed-forward layers are 662
key-value memories. In *Proceedings of the 2021 663
Conference on Empirical Methods in Natural Lan- 664
guage Processing*, pages 5484–5495. 665

Ali Gholipour, Nasser Kehtarnavaz, Richard Briggs, 666
Michael Devous, and Kaundinya Gopinath. 2007. 667
Brain functional localization: a survey of image reg- 668
istration techniques. *IEEE transactions on medical 669
imaging*, 26(4):427–451. 670

Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, 671
Mantas Mazeika, Dawn Song, and Jacob Steinhardt. 672
2020. Measuring massive multitask language under- 673
standing. In *International Conference on Learning 674
Representations*. 675

676	Ganesh Jawahar, Benoît Sagot, and Djamé Seddah.	<i>on Natural Language Processing (EMNLP-IJCNLP)</i> ,	732
677	2019. What does bert learn about the structure of	pages 2463–2473.	733
678	language? In <i>Proceedings of the 57th Annual Meet-</i>		
679	<i>ing of the Association for Computational Linguistics</i> ,	Pouya Pezeshkpour. 2023. Measuring and modifying	734
680	pages 3651–3657.	factual knowledge in large language models. <i>arXiv</i>	735
		<i>preprint arXiv:2306.06264</i> .	736
681	Albert Q Jiang, Alexandre Sablayrolles, Arthur Men-		
682	sch, Chris Bamford, Devendra Singh Chaplot, Diego	Alec Radford, Jeffrey Wu, Rewon Child, David Luan,	737
683	de las Casas, Florian Bressand, Gianna Lengyel, Guil-	Dario Amodei, Ilya Sutskever, et al. Language mod-	738
684	laume Lample, Lucile Saulnier, et al. 2023. Mistral	els are unsupervised multitask learners.	739
685	7b. <i>arXiv preprint arXiv:2310.06825</i> .		
686	Zhengbao Jiang, Antonios Anastasopoulos, Jun Araki,	Siyu Ren and Kenny Zhu. 2022. Specializing pre-	740
687	Haibo Ding, and Graham Neubig. 2020. X-factr:	trained language models for better relational reason-	741
688	Multilingual factual knowledge retrieval from pre-	ing via network pruning. In <i>Findings of the Associ-</i>	742
689	trained language models. In <i>Proceedings of the 2020</i>	<i>ation for Computational Linguistics: NAACL 2022</i> ,	743
690	<i>Conference on Empirical Methods in Natural Lan-</i>	pages 2195–2207.	744
691	<i>guage Processing (EMNLP)</i> , pages 5943–5959.		
692	Yiming Ju and Zheng Zhang. 2023. Klob: a bench-	Adam Roberts, Colin Raffel, and Noam Shazeer. 2020.	745
693	mark for assessing knowledge locating methods in	How much knowledge can you pack into the param-	746
694	language models. <i>arXiv preprint arXiv:2309.16535</i> .	eters of a language model? In <i>Proceedings of the</i>	747
		<i>2020 Conference on Empirical Methods in Natural</i>	748
		<i>Language Processing (EMNLP)</i> , pages 5418–5426.	749
695	Nora Kassner, Philipp Dufter, and Hinrich Schütze.		
696	2021. Multilingual lama: Investigating knowledge in	Gerard Salton. 1983. Introduction to modern informa-	750
697	multilingual pretrained language models. In <i>Proceed-</i>	tion retrieval. <i>McGraw-Hill</i> .	751
698	<i>ings of the 16th Conference of the European Chap-</i>		
699	<i>ter of the Association for Computational Linguistics:</i>	Mukund Sundararajan, Ankur Taly, and Qiqi Yan. 2017.	752
700	<i>Main Volume</i> , pages 3250–3258.	Axiomatic attribution for deep networks. In <i>Internat-</i>	753
		<i>ional conference on machine learning</i> , pages 3319–	754
701	Jacob Devlin Ming-Wei Chang Kenton and Lee Kristina	3328. PMLR.	755
702	Toutanova. 2019. Bert: Pre-training of deep bidirec-		
703	tional transformers for language understanding. In	Mujeen Sung, Jinhyuk Lee, Sean Yi, Minji Jeon, Sung-	756
704	<i>Proceedings of NAACL-HLT</i> , pages 4171–4186.	dong Kim, and Jaewoo Kang. 2021. Can language	757
		models be biomedical knowledge bases? In <i>Proceed-</i>	758
705	Nelson F Liu, Matt Gardner, Yonatan Belinkov,	<i>ings of the 2021 Conference on Empirical Methods</i>	759
706	Matthew E Peters, and Noah A Smith. 2019. Linguis-	<i>in Natural Language Processing</i> , pages 4723–4734.	760
707	tic knowledge and transferability of contextual rep-		
708	resentations. In <i>Proceedings of the 2019 Conference of</i>	Gemma Team, Thomas Mesnard, Cassidy Hardin,	761
709	<i>the North American Chapter of the Association for</i>	Robert Dadashi, Surya Bhupatiraju, Shreya Pathak,	762
710	<i>Computational Linguistics: Human Language Tech-</i>	Laurent Sifre, Morgane Rivière, Mihir Sanjay Kale,	763
711	<i>nologies, Volume 1 (Long and Short Papers)</i> , pages	Juliette Love, et al. 2024. Gemma: Open models	764
712	1073–1094.	based on gemini research and technology. <i>arXiv</i>	765
		<i>preprint arXiv:2403.08295</i> .	766
713	Leland McInnes, John Healy, Nathaniel Saul, and Lukas	Hugo Touvron, Louis Martin, Kevin Stone, Peter Al-	767
714	Großberger. 2018. UMAP: uniform manifold ap-	bert, Amjad Almahairi, Yasmine Babaei, Nikolay	768
715	proximation and projection . <i>J. Open Source Softw.</i> ,	Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti	769
716	3(29):861.	Bhosale, et al. 2023. Llama 2: Open founda-	770
		tion and fine-tuned chat models. <i>arXiv preprint</i>	771
717	Kevin Meng, David Bau, Alex Andonian, and Yonatan	<i>arXiv:2307.09288</i> .	772
718	Belinkov. 2022a. Locating and editing factual as-		
719	sociations in gpt. <i>Advances in Neural Information</i>	Ningyu Zhang, Yunzhi Yao, Bozhong Tian, Peng Wang,	773
720	<i>Processing Systems</i> , 35:17359–17372.	Shumin Deng, Mengru Wang, Zekun Xi, Shengyu	774
		Mao, Jintian Zhang, Yuansheng Ni, et al. 2024a. A	775
721	Kevin Meng, Arnab Sen Sharma, Alex J Andonian,	comprehensive study of knowledge editing for large	776
722	Yonatan Belinkov, and David Bau. 2022b. Mass-	language models. <i>arXiv preprint arXiv:2401.01286</i> .	777
723	editing memory in a transformer. In <i>The Eleventh</i>		
724	<i>International Conference on Learning Representa-</i>	Zhengyan Zhang, Zhiyuan Zeng, Yankai Lin, Chaojun	778
725	<i>tions</i> .	Xiao, Xiaozhi Wang, Xu Han, Zhiyuan Liu, Ruob-	779
		ing Xie, Maosong Sun, and Jie Zhou. 2023. Emer-	780
726	Fabio Petroni, Tim Rocktäschel, Sebastian Riedel,	gent modularity in pre-trained transformers. In <i>Find-</i>	781
727	Patrick Lewis, Anton Bakhtin, Yuxiang Wu, and	<i>ings of the Association for Computational Linguistics:</i>	782
728	Alexander Miller. 2019. Language models as knowl-	<i>ACL 2023</i> , pages 4066–4083.	783
729	edge bases? In <i>Proceedings of the 2019 Confer-</i>		
730	<i>ence on Empirical Methods in Natural Language Pro-</i>		
731	<i>cessing and the 9th International Joint Conference</i>		

784 Zhihao Zhang, Jun Zhao, Qi Zhang, Tao Gui, and
785 Xuanjing Huang. 2024b. Unveiling linguistic re-
786 gions in large language models. *arXiv preprint*
787 *arXiv:2402.14700*.

788 Wayne Xin Zhao, Kun Zhou, Junyi Li, Tianyi Tang,
789 Xiaolei Wang, Yupeng Hou, Yingqian Min, Beichen
790 Zhang, Junjie Zhang, Zican Dong, et al. 2023. A
791 survey of large language models. *arXiv preprint*
792 *arXiv:2303.18223*.

793 Xuhui Zhou, Yue Zhang, Leyang Cui, and Dandan
794 Huang. 2020. Evaluating commonsense in pre-
795 trained language models. In *Proceedings of the*
796 *AAAI conference on artificial intelligence*, volume 34,
797 pages 9733–9740.

798 A Latent Value Analysis

799 As previously found by (Geva et al., 2021), trans-
800 former feed-forward layers can be viewed as key-
801 value memory units, with hidden activations acting
802 as coefficients for the individual memories stored
803 in \mathbf{W}^D . Thus, a natural question to explore is what
804 are the properties of the memory cells associated
805 with the key neurons for the different domains and
806 if they are clustered in the corresponding semantic
807 space.

808 As a first step in our analysis, we visualise
809 the \mathbf{W}^D vectors associated with the knowledge
810 neurons from the different domains using UMAP
811 (McInnes et al., 2018) for dimensionality reduction
812 (with cosine similarity used as the distance metric).
813 For comparison, we additionally include the vec-
814 tors from the unembedding matrix. The results are
815 shown in Figure A1. As can be seen from the fig-
816 ure, the distribution of the vectors associated with
817 key knowledge neurons appears to be significantly
818 different from that of vector unembeddings. Thus,
819 it appears that the contents of the internal memory
820 cells used by LLaMA 2 are not directly aligned
821 with the candidate output tokens.

822 Since the 2D visualisation produced by UMAP
823 might not accurately reflect the true properties of
824 the data manifold, we additionally examined the
825 highest-likelihood tokens for the key domain neu-
826 ron memory cells. These were computed by di-
827 rectly applying the Llama 2 unembedding layer to
828 the vectors stored in these cells. We found the re-
829 sulting tokens rather uninterpretable, including to-
830 kens like `textt`, `archivi`, `_kontrola`, `_totalité`
831 or `_Einzeln`. Upon closer investigation, we found
832 these to be closely matching the set of unembed-
833 ding vectors with the largest vector norms (which
834 we would expect to generally receive higher like-
835 lihoods when multiplied with vectors not aligned

836 with any of the unembeddings). This seems to
837 further support a conjecture that the memory cell
838 vectors associated with the located domain-specific
839 neurons might capture intermediate data in a sub-
840 space different from the one used for the final token
841 prediction. Apart from the above tokens, we also
842 found option letters A, B, C and D to be represented
843 in the highest-likelihood tokens. This suggests that
844 some neurons within the identified key neuron set
845 may still correspond to option letters, as mentioned
846 in the Limitations section.

847 We leave further investigation and confirmation
848 of these findings for future work.

849 B Neuron-Based Prediction Details

850 In the neuron-based prediction case study, we ex-
851 periment on the MMLU (Hendrycks et al., 2020)
852 validation set to ensure there is no overlap between
853 the dataset used to mine the key neurons and the
854 test set. Thus, the considered domain neurons were
855 determined based on queries not used for this ex-
856 periment. As a further post-processing step, we
857 randomly select three options from other domains
858 to replace the incorrect options in each query. Ad-
859 ditionally, we manually remove questions that be-
860 come invalid due to this post-processing, including
861 queries such as “Which of the following is LEAST
862 valid?” and “All of the following statements are
863 true EXCEPT”. These operations result in ~ 20
864 test samples per domain. To perform the neuron-
865 based prediction, we compute the gradient of the
866 probability of each option token with respect to
867 the key neurons for the domain of the considered
868 query, and select the option with the highest total
869 gradient.

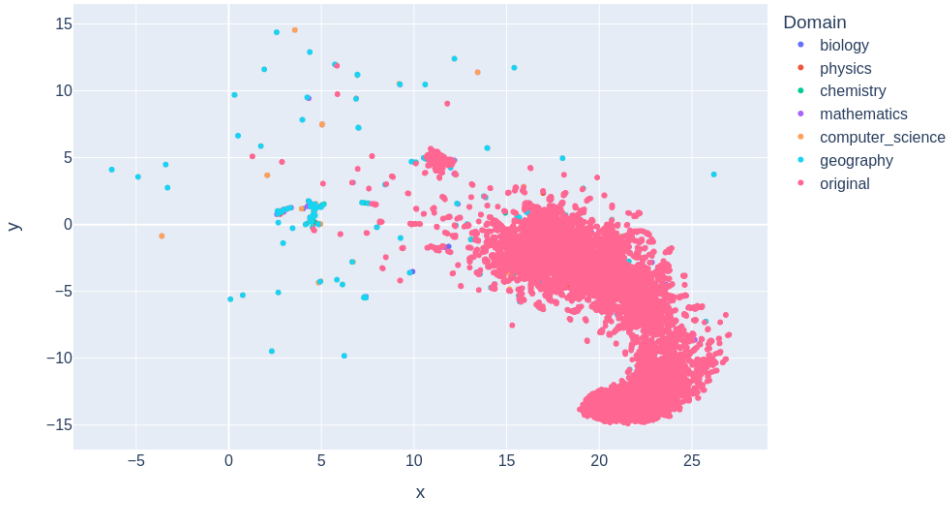


Figure A1: UMAP visualisation of W^D vectors associated with the knowledge neurons and the token unembeddings

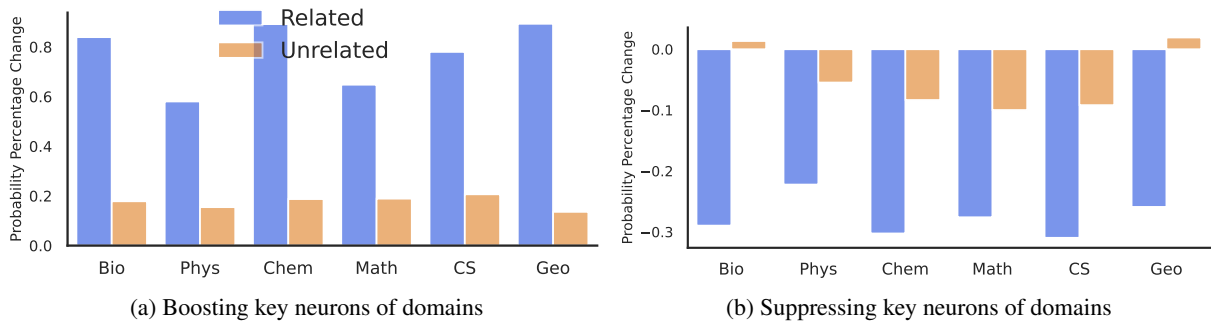


Figure A2: The correct probability percentage change across different domains. The LLM here is Mistral-7B (Jiang et al., 2023)

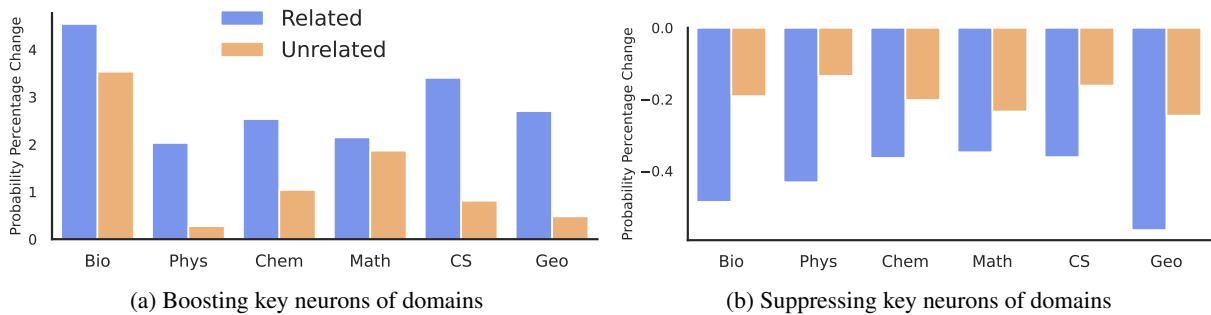


Figure A3: A comparison study of using $\frac{\partial P_q}{\partial g}$ to compute NA-ICA scores. The LLM here is LLaMA-7B (Touvron et al., 2023).

Num	Template
Domain Prompt 1	<i>You will be asked a multiple-choice question. Respond with the letter which corresponds to the correct answer, followed by a period. There is no need to provide an explanation, so your response should be very short.</i> <i>Now here is the question:</i> {Question} A. {A} B. {B} C. {C} D. {D} Response:
Domain Prompt 2	<i>Prepare to answer a multiple-choice question. Provide the letter that corresponds to the correct answer, followed by a period. Keep your response brief; no explanations are necessary.</i> <i>Here is the question:</i> {Question} A. {A} B. {B} C. {C} D. {D} Response:
Domain Prompt 3	<i>Below is a multiple-choice question. Respond with the letter that best answers the question. Keep your response brief, stating only the letter corresponding to your answer, followed by a period, with no explanation.</i> <i>The question is:</i> {Question} A. {A} B. {B} C. {C} D. {D} Response:
Language Prompt 1	<i>You will be asked a multiple-choice question. Respond with the letter which corresponds to the correct answer, followed by a period. There is no need to provide an explanation, so your response should be very short.</i> <i>Now here is the question:</i> {Question} <i>Here the [Y] is most likely to be?</i> A. {A} B. {B} C. {C} D. {D} Response:
Language Prompt 2	<i>Prepare to answer a multiple-choice question. Provide the letter that corresponds to the correct answer, followed by a period. Keep your response brief; no explanations are necessary.</i> <i>Now here is the question:</i> {Question} <i>Here the [Y] is most likely to be?</i> A. {A} B. {B} C. {C} D. {D} Response:
Language Prompt 3	<i>Below is a multiple-choice question. Respond with the letter that best answers the question. Keep your response brief, stating only the letter corresponding to your answer, followed by a period, with no explanation.</i> <i>Now here is the question:</i> {Question} <i>Here the [Y] is most likely to be?</i> A. {A} B. {B} C. {C} D. {D} Response:

Table A1: Prompt templates for constructing multi-choice QA datasets. We use ChatGPT to translate English templates to other languages.

Field	Question	Options
Biology	<i>The energy given up by electrons as they move through the electron transport chain is used to?</i>	A. make glucose B. make NADH C. produce ATP D. break down glucose
Physics	<i>An object is placed 100 cm from a plane mirror. How far is the image from the object?</i>	A. 50 cm B. 200 cm C. 100 cm D. 300 cm
Chemistry	<i>Three half-lives after an isotope is prepared:</i>	A. 12.5% of the isotope decayed B. 25% of the isotope decayed C. 25% of the isotope is left D. 12.5% of the isotope is left
Mathematics	<i>Suppose the graph of f is both increasing and concave up on $a <= x <= b$. Then, using the same number of subdivisions, and with L, R, M, and T denoting, respectively, left, right, midpoint, and trapezoid sums, it follows that:</i>	A. $R <= T <= M <= L$ B. $L <= M <= T <= R$ C. $R <= M <= T <= L$ D. $L <= T <= M <= R$
Computer Science	<i>A programmer is writing a program that is intended to be able to process large amounts of data. Which of the following considerations is LEAST likely to affect the ability of the program to process larger data sets?</i>	A. How long the program takes to run B. How many programming statements the program contains C. How much storage space the program requires as it runs D. How much memory the program requires as it runs
Geography	<i>The tendency for migration to decrease with distance is called?</i>	A. push factors. B. migration selectivity. C. distance decay. D. pull factors.
English	<i>Sergey Lavrov was born in [Y]. Here the [Y] is most likely to be?</i>	A. Montevideo B. Bengaluru C. Parsons D. Moscow

Table A2: Examples in our constructed datasets. For the language dataset, we only show one English example as multilingual samples are obtained by using translator (Kassner et al., 2021)