

# Cracking Factual Knowledge: A Comprehensive Analysis of Degenerate Knowledge Neurons in Large Language Models

Anonymous ACL submission

## Abstract

Large language models (LLMs) store extensive factual knowledge, but the underlying mechanisms remain unclear. Previous research suggests that factual knowledge is stored within multi-layer perceptron weights, and some storage units exhibit degeneracy, referred to as *Degenerate Knowledge Neurons* (DKNs). Despite the novelty and unique properties of this concept, it has not been rigorously defined or systematically studied. We first consider the connection weight patterns of MLP neurons and define DKNs from both structural and functional aspects. Based on this, we introduce the *Neurological Topology Clustering* method, which allows the formation of DKNs in any numbers and structures, leading to a more accurate DKN acquisition. Furthermore, inspired by cognitive science, we explore the relationship between DKNs and the robustness, evolvability, and complexity of LLMs. Our execution of 34 experiments under 6 settings demonstrates the connection between DKNs and these three properties. The code will be available soon.

## 1 Introduction

Large pretrained language models (PLMs) are believed to store extensive factual knowledge (Touvron et al., 2023; OpenAI et al., 2023), yet the mechanisms of knowledge storage in PLMs remain largely unexplored. Dai et al. (2022) propose that some multi-layer perceptron (MLP) neurons can store “knowledge”. As shown in the Figure 1(a), for the fact  $\langle \text{COVID-19}, \text{dominant variant}, \text{Delta} \rangle$ , the corresponding knowledge storage units  $a$  through  $f$  are termed knowledge neurons (KNs). Chen et al. (2024a) find that distinct KN pairs, such as  $\{a, b\}$  and  $\{c, d\}$  in Figure 1(b<sub>1</sub>), can store identical facts. They define the set of these pairs as degenerate knowledge neurons (DKNs) from a functional perspective.

While Chen et al. (2024a) have conducted some exploration on DKNs, their definition and acqui-

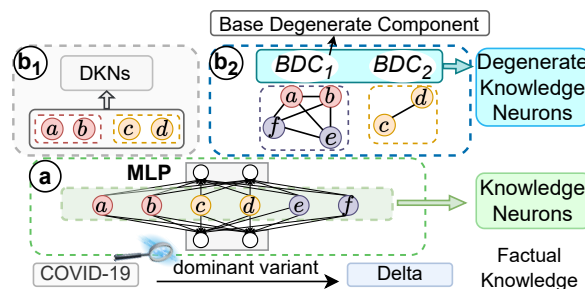


Figure 1: Explanation of KNs and DKNs. (a) illustrates the KNs corresponding to a fact, (b<sub>1</sub>) and (b<sub>2</sub>) are the preliminary definition of DKNs and our complete definition of DKNs, respectively.

sition method for DKNs still face two issues. (1) *Numerical Limitation*: They constrain each DKN’s element to contain just two KNs, such as  $\{a, b\}$  in Figure 1(b<sub>1</sub>). However, factual knowledge may require more than two neurons for representation (Allen-Zhu and Li, 2023). (2) *Connectivity Oversight*: They only consider the neurons themselves, neglecting the connectivity weights between neurons. However, knowledge expression requires the interaction of multiple neurons (Zhu and Li, 2023), thus, it is necessary to consider the connectivity structure between neurons.

To address these two issues, we first provide a comprehensive definition of DKNs from two perspectives. **Functionally**, some subsets of KNs can independently express the same fact, termed as Base Degenerate Components (BDCs), such as  $BDC-1$  and  $BDC-2$  in Figure 1(b<sub>2</sub>). In other words, they exhibit mutual degeneracy. The set of these BDCs is referred to as a DKN. **Structurally**, as shown in Figure 1(b<sub>2</sub>), BDCs like  $BDC-1$  and  $BDC-2$  differ in KN number and connection tightness. To assess DKNs’ structural traits, we define neuron distances based on connection weights and analyze the structural properties of neuron sets accordingly.

Based on the above definition, we introduce the *Neurological Topology Clustering* (NTC, §3) method to obtain DKNs, which includes cluster-

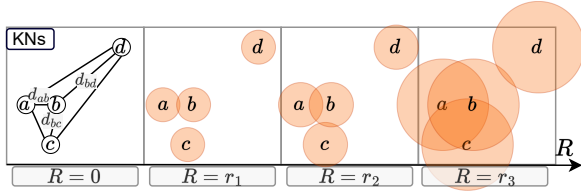


Figure 2: The clustering part of the Neurological Topology Clustering method. The  $x$ -axis ( $R$ ) represents the increasing distance threshold starting from 0. Circles with radius  $R$  are drawn around neurons, and intersecting circles indicate that the KNs are clustered together.

ing and filtering stages, enabling the formation of BDCs with an arbitrary number of neurons and arbitrary neuronal connection structures. Figure 2 illustrates the clustering stage. Given four KNs  $\{a, b, c, d\}$  with fixed connection weights and an increasing distance threshold  $R$  starting from 0, we observe whether the KNs can cluster together as  $R$  changes. At  $R = 0$ , the KNs are isolated points. When  $R = r_1 > d_{ab}$ ,  $\{a, b\}$  form a cluster; at  $R = r_2 > d_{bc}$ ,  $\{a, b, c\}$  cluster together; and at  $R = r_3 > d_{bd}$ ,  $\{a, b, c, d\}$  form a single cluster. Notably, a wide range of  $R$  values maintains the  $\{a, b, c\}$  cluster (from  $r_2$  to  $r_3$ ), indicating its stable existence. This stable cluster, suggesting a strong knowledge expression ability (Zhu and Li, 2023), is identified as a BDC. BDCs are then filtered to derive DKNs, as detailed in Section 3.

Furthermore, inspired by cognitive science research on degeneracy (Whitacre and Bender, 2010; Edelman and Gally, 2001; Whitacre, 2010; Mason, 2015), we investigate the relationship between DKNs and the robustness, evolvability, and complexity of PLMs. Our findings are illustrated in Figure 3 and elaborated upon in the following:

(1) **Robustness (§4)**: One aspect of PLMs’ robustness is their ability to handle input interference (Fernandez et al., 2005). Given a query subject to interference, we suppress (or enhance) the values or connection weights of DKNs and observe resulting changes in the PLMs’ prediction probability for the query. Furthermore, we conduct a fact-checking experiment (Guo et al., 2022), using DKNs to detect false facts. Experiments demonstrate that **DKNs can help PLMs cope with input interference**, indicating their contribution to robustness.

(2) **Evolvability (§5)**: Evolvability is defined as the ability to adaptively evolve in new environments (Kirschner and Gerhart, 1998). Inspired by this, we consider learning new knowledge as one aspect of PLMs’ evolvability. We validate the sig-

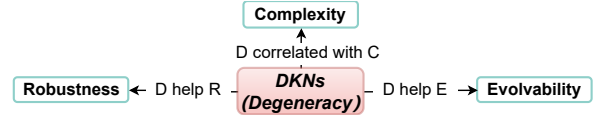


Figure 3: The relationship between DKNs (with degeneracy property), robustness, evolvability and complexity in PLMs.

nificance of DKNs in PLMs’ ability to learn new knowledge from two perspectives. First, to prove that PLMs utilize DKNs to learn new knowledge, we directly fine-tune the PLMs and find that the regions of parameter changes highly overlap with DKNs. Second, we employ an efficient fine-tuning method, freezing all MLP neurons except DKNs. We discover that PLMs can utilize DKNs to efficiently learn new knowledge while not forgetting old knowledge. In summary, **DKNs enable PLMs to learn new knowledge more efficiently**, highlighting their crucial role in enhancing evolvability.

(3) **Complexity (§6)**: PLMs’ complexity is positively correlated to the number of parameters (Mars, 2022). In experiments related to the robustness and evolvability of PLMs, we compare the performance of PLMs across different scales and find that **degeneracy is positively correlated with complexity**. In addition to these main experiments, we also conduct a supplementary fact-checking experiment on complex texts. This additional study proves that we can combine the large PLMs’ complex text understanding ability with the DKNs’ fact-checking capability to complete the task.

Our contributions can be summarized as follows:

- We provide a comprehensive definition of DKNs from both functional and structural aspects, pioneering the study of structures in PLMs’ factual knowledge storage units.
- We introduce the Neurological Topology Clustering method, allowing the formation of DKNs in any numbers and structures, leading to a more accurate DKN acquisition.
- Through the lens of DKNs, we investigate the robustness, evolvability, and complexity of PLMs. By conducting 34 experiments under 6 settings, we demonstrate the relationship between DKNs and these three properties.

## 2 Datasets and Models

We utilize the TempLama dataset (Dhingra et al., 2022) to analyze DKNs. Each data instance includes a relation name, a date, a query, and an answer, such as  $\langle P37, \text{September 2021}, \text{COVID-19},$

*dominant variant*,  $\_$ ). Except for timestamps, our dataset matches the Lama (Petroni et al., 2019a, 2020) and mLama (Kassner et al., 2021) format used by Dai et al. (2022) and Chen et al. (2024a). Under different experimental settings, we perform data augmentation or filtering operations on the dataset, resulting in six datasets that differ in content but are consistent in format. Regarding model selection, we choose GPT-2 (Radford et al., 2019) and LLaMA2-7b (Touvron et al., 2023), both of which are based on the currently most popular autoregressive architecture but have different parameter sizes, allowing us to test the generalization and scalability of our methods and conclusions.

### 3 Neurological Topology Clustering

#### 3.1 Definition of DKNs

**Formalization** Given a fact, we utilize the AMIG method (Chen et al., 2024a) to obtain KNs, denoting them as  $\mathcal{N} = \{n_1, n_2, \dots, n_k\}$ , where  $n_i$  is a KN. For details of this method, see Appendix B. Let DKNs be denoted as  $\mathcal{D}$ , containing  $s$  elements,  $\mathcal{D} = \{\mathcal{B}_1, \mathcal{B}_2, \dots, \mathcal{B}_s\}$ , where  $\mathcal{B}_j = \{n_{j1}, n_{j2}, \dots, n_{|\mathcal{B}_j|}\}$  is named as the Base Degenerate Component (BDC). Thus, this fact ultimately corresponds to a set of DKNs:

$$\mathcal{D} = \{\mathcal{B}_1, \dots, \mathcal{B}_s\} = \{(n_{11}, \dots, n_{|\mathcal{B}_1|}), \dots, (n_{s1}, \dots, n_{|\mathcal{B}_s|})\} \quad (1)$$

**Functional Definition** Degeneracy requires that each BDC should independently express a fact. Let  $Prob(\mathcal{B})$  represent the PLMs’ answer prediction probability when  $\mathcal{B}$  is activated, then the functional definition of DKNs is:

$$Prob(\mathcal{D}) \approx Prob(\mathcal{B}_i), \forall i = 1, 2, \dots, s \quad (2)$$

$$Prob(\emptyset) \ll Prob(\mathcal{B}_i), \forall i = 1, 2, \dots, s \quad (3)$$

where Equation 2 indicates that activating any single BDC is sufficient to express the fact, and Equation 3 suggests that if all BDCs are suppressed (i.e., activating the empty set  $\emptyset$ ), the fact cannot be correctly expressed.

**Structural Definition** Zhu and Li(2023) argue that tightly connected neurons tend to store knowledge centrally. Thus, we use the adjacency matrix  $\mathcal{A}$  to evaluate the DKNs’ connection structure. For neurons  $A$  and  $B$  in layer  $l_A$  and layer  $l_B$  respectively, we calculate the distance  $d_{AB}$  as follows:

$$d_{AB} = \begin{cases} |1/w_{AB}| & \text{if } w_{AB} \neq 0 \text{ and } |l_A - l_B| = 1, \\ \min_{P \in \text{Paths}(\mathcal{N})} \sum_{(i,j) \in P} d_{ij} & \text{if } |l_A - l_B| > 1 \text{ and a path exists,} \\ \infty & \text{otherwise.} \end{cases} \quad (4)$$

---

#### Algorithm 1: Neurological Topology Clustering

---

**Input:** Knowledge neurons  $\mathcal{N}$ , Adjacent matrix  $\mathcal{A}$ , dynamic threshold  $\tau_1$  and threshold  $\tau_2$   
**Output:** Degenerate knowledge neurons  $\mathcal{D}$   
*// Initialization*  
Initialize  $\mathcal{D} = \emptyset$  and  $R \leftarrow 0$ , where  $R$  is the distance threshold for persistent homology.  
*// Persistent Homology (Clustering)*  
**while**  $R$  increases **do**  
    Record all base degenerate components  $\mathcal{B}_i$  and their corresponding persistence duration  $R_p$   
**end**  
*// Filtering*  
**for each**  $\mathcal{B}_i$  **do**  
    **if**  $R_p(\mathcal{B}_i) > \tau_1$  and  $Prob(\mathcal{B}_i) > \tau_2$  **then**  
        Add  $\mathcal{B}_i$  to  $\mathcal{D}$ , where  $Prob(\mathcal{B}_i)$  is the prediction probability when  $\mathcal{B}_i$  is activated  
    **end**  
**end**  
**return**  $\mathcal{D}$

---

where  $\text{Paths}(\mathcal{N})$  includes all paths from  $A$  to  $B$  through  $\mathcal{N}$ . Distance calculation varies in three scenarios. First, for neurons in adjacent layers, it is the reciprocal of the weight. Second, for neurons spanning multiple layers, it is the shortest distance determined by a dynamic programming algorithm. Third, for neurons in the same layer, since information in PLMs transmits between layers rather than within a layer (Meng et al., 2022), we set the distance to  $\infty$ . Hence, any  $\mathcal{D}$  can correspond to an adjacency matrix  $\mathcal{A}$ , where  $\mathcal{A} \in \mathbb{R}^{k \times k}$ , and  $k$  is the number of knowledge neurons contained in  $\mathcal{D}$ . Based on  $\mathcal{A}$ , beyond traditional neuron value editing, modifying the connection weights offers another method to edit PLMs.

#### 3.2 The Acquisition of DKNs

**Persistent Homology** In our method, we capture the persistent homology (Edelsbrunner et al., 2008) of KN sets, which represents the duration of the set’s existence and the tightness of the set’s connections. Figure 2 informally illustrates this process. Formally, given two KNs,  $n_i$  and  $n_j$ , which are both the centers of expanding circles. When the start radius (i.e., distance threshold) is 0, this corresponds to  $R_s = 0$ . Suppose they touch at radius  $R_e = r_1$ , which marks the end radius  $R = r_1$ , and a new start radius  $R_s = r_1$ . At this point,  $n_i$  and  $n_j$  are clustered together, forming a BDC,  $\mathcal{B} = \{n_i, n_j\}$ , corresponding to a persistence duration  $R_p = R_e - R_s = r_1$ . For details on persistent homology, see Appendix C.

**NTC Method** Our method is shown in Figure 2 and Algorithm 1. We design two steps, i.e., clustering and filtering steps, to obtain DKNs. During the

Method	GPT-2						
	2	3	4	5	6	7	Average
DBSCAN	12.0 → 26.0	17.7 → 38.8	23.4 → 33.8	32.9 → 53.4	24.8 → 49.8	26.4 → 46.6	23.90 → 34.72
Hierarchical	2.9 → 3.8	15.5 → 7.3	3.7 → 12	4.3 → 27	25.0 → 92.0	—	20.78 → 30.81
K-Means	27.4 → 38.3	17.8 → 29.5	27.8 → 44.6	32.4 → 58.7	28.9 → 30.5	21.8 → 30.1	34.42 → 38.86
AMIG	-0.8 → 0.9	—	—	—	—	—	-0.8 → 0.9
<b>NTC (Ours)</b>	<b>7.9 → 53</b>	<b>8.6 → 44</b>	<b>12 → 92</b>	<b>11 → 53</b>	<b>3.1 → 32</b>	<b>7.8 → 61</b>	<b>9.32 → 55.60</b>

Method	LLaMA2						
	2	3	8	11	14	17	Average
DBSCAN	7.1 → 15.4	8.7 → 15.2	7.7 → 16.3	7.2 → 25.3	—	—	22.26 → 18.24
Hierarchical	27.6 → 44.3	22.3 → 6.5	—	—	—	—	27.50 → 44.04
K-Means	2.8 → 16.1	19.2 → 39.1	37.9 → 97.1	50.7 → 111	36.3 → 50.0	4.9 → 17.9	20.08 → 39.24
AMIG	-1.0 → 2.0	—	—	—	—	—	-1.0 → 2.0
<b>NTC (Ours)</b>	<b>2.8 → 15.6</b>	<b>4.3 → 19.1</b>	<b>7.8 → 50.1</b>	<b>4.6 → 25.7</b>	<b>13.6 → 37.7</b>	<b>31.2 → 135</b>	<b>14.11 → 28.78</b>

Table 1: Comparison of  $\mathcal{D}$  obtained by different methods. The numbers above each column indicate the cardinality of the BDC sets. Each cell shows two values for the PLMs’ prediction probability decrease ( $\Delta Prob$  (%)): the left is the average  $\Delta Prob$  when partially suppressing BDCs (i.e., 1 to n-1 BDCs), and the right is the  $\Delta Prob$  when fully suppressing all BDCs. Lower left values, higher right values, and larger differences between them indicate better degeneracy. The symbol “—” signifies that a specific method failed to generate a  $\mathcal{D}$  with the specified cardinality.

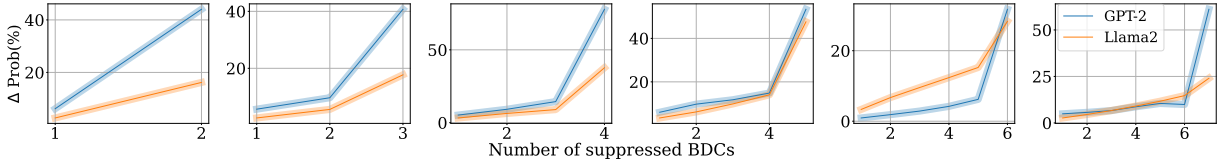


Figure 4: The relationship between  $\Delta Prob$  and the number of suppressed BDCs (using NTC). Table 1 averages the results for suppressing 1 to n-1 BDCs, while this figure shows the changes as the number of suppressed BDCs varies from 1 to n, with the final point representing all BDCs suppressed. Lower  $\Delta Prob$  for partial suppression and higher  $\Delta Prob$  for full suppression, i.e., a more prominent final turning point, indicate better degeneracy. As the cardinality of  $\mathcal{D}$  varies across PLMs and methods, Table 1 and Figure 4 show representative results. Full results are in Appendix D (Table 4, Figures 12, 13, 14, 15, 16).

clustering process, as  $R$  increases from 0 to infinity, we record all BDCs along with their corresponding  $R_p$ . During the filtering process, we initially select BDCs with  $R_p$  above  $\tau_1$ . Then, among these, only BDCs with a  $Prob(\mathcal{B}_i)$  greater than a threshold  $\tau_2$  are kept. Finally, these BDCs constitute  $\mathcal{D}$ , as shown in the formula below.

$$\mathcal{D} = \{\mathcal{B}_i | R_p(\mathcal{B}_i) > \tau_1 \text{ and } Prob(\mathcal{B}_i) \geq \tau_2\} \quad (5)$$

### 3.3 Experiments of DKNs Acquisition

**Experimental settings** Given a set  $\mathcal{D} = \{\mathcal{B}_1, \mathcal{B}_2, \dots, \mathcal{B}_n\}$ , we traverse all subsets of  $\mathcal{D}$ , suppressing values or connection weights of all BDCs in each subset. We then calculate the drop in the answer probability before (b) and after (a) suppression:  $\Delta Prob(\%) = \frac{Prob_b - Prob_a}{Prob_b}$ . We select four other methods as baselines, including K-Means (Ahmed et al., 2020), DBSCAN (Ester et al., 1996), Hierarchical Clustering (Murtagh and Contreras, 2012) and AMIG (Chen et al., 2024a).

**Findings** (1) DKNs identified by NTC exhibit strong degeneracy. As shown in Table 1, when sup-

pressing a subset of BDCs, NTC yields a smaller  $\Delta Prob$  compared to other baseline methods. Conversely, when suppressing all BDCs, NTC results in the largest  $\Delta Prob$ . The clear inflection point in Figure 4 further demonstrates that suppressing a subset of BDCs does not lead to a gradual increase in  $\Delta Prob$ . Instead, as long as one BDC exists, PLMs have a sufficient ability to express facts. In summary, the DKNs obtained by the NTC method most closely align with the definitions of DKNs in Equations 2 and 3, exhibiting the most desirable degeneracy properties. (2) Suppressing DKNs may inadvertently enhance PLMs’ knowledge expression. Negative values in Table 1 indicate increased prediction probability, suggesting that suppressing some DKN subsets allows others to compensate and improve PLM performance.

## 4 The Impact of DKNs on Robustness

### 4.1 Query-Perturbation

**Explanation** Robustness to input errors is one facet of PLMs’ overall robustness. This is critical

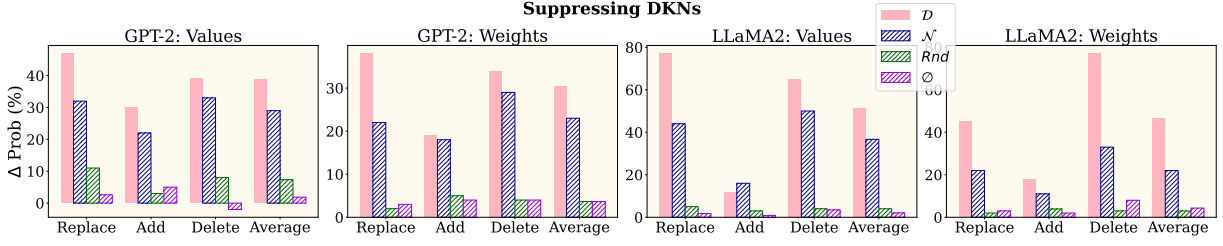


Figure 5: Changes in the prediction probabilities of PLMs corresponding to the suppression of DKNs and other baselines.  $D$ ,  $N$ ,  $Rnd$ , and  $\emptyset$  represent DKNs, KNs, random neurons, and suppressing no neurons, respectively.

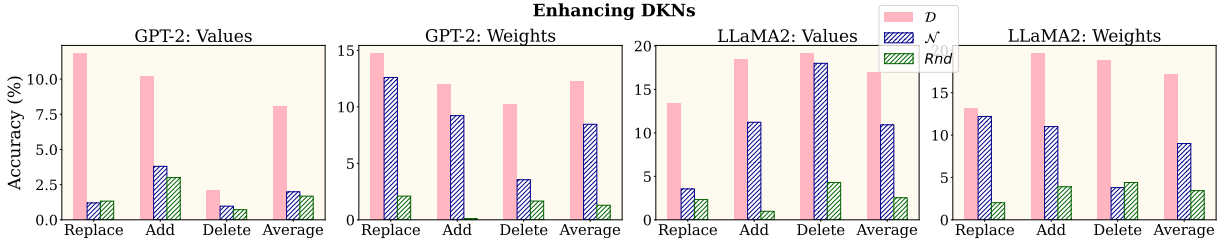


Figure 6: Changes in the accuracy of PLMs corresponding to the enhancement of DKNs and other baselines.

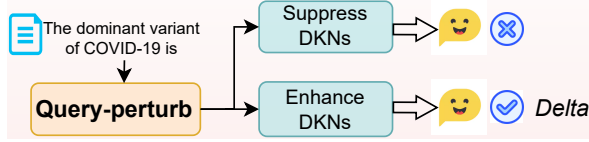


Figure 7: Query perturbation and the impact of suppressing or enhancing DKNs on PLMs’ performance.

in practical scenarios where PLMs often encounter user input errors like spelling mistakes or character omissions (Chen et al., 2010). To simulate this scenario, we apply random disturbances to the inputs. Given an input sequence  $Q = \{q_1, q_2, \dots, q_n\}$ , we generate its perturbed counterpart  $Q^*$ :

$$Q^* = \begin{cases} \{q_1, \dots, q_{i-1}, [\text{replace}], q_{i+1}, \dots, q_n\} & \text{if replace,} \\ \{q_1, \dots, q_{i-1}, [\text{add}], q_i, \dots, q_n\} & \text{if add,} \\ \{q_1, \dots, q_{i-1}, q_{i+1}, \dots, q_n\} & \text{if delete.} \end{cases} \quad (6)$$

where “[replace]” and “[add]” are special characters. To investigate the role of DKNs, we suppress and enhance them separately (Figure 7).

**Suppressing DKNs** We apply two different operations to suppress DKNs: (1) zeroing neuron values, and (2) nullifying neuron connection weights. For comparison, we select three baselines: suppressing knowledge neurons, suppressing random neurons (the same number as DKNs), and suppressing no neurons. For the PLMs subjected to suppression operations, we calculate their prediction probabilities for both  $Q$  and  $Q^*$ , and compute the change in probability:  $\Delta Prob(\%) = \frac{Prob(Q) - Prob(Q^*)}{Prob(Q)}$ . Figure 5 presents our results.

**Enhancing DKNs** Since robustness in PLMs is always limited, they may answer some  $Q^*$  queries

incorrectly. We record these queries as  $Q_{err}^*$ , and apply two operations to enhance DKNs: (1) doubling neuron values, and (2) doubling neuron connection weights. Then, we reassess the enhanced PLMs’ accuracy on  $Q_{err}^*$ , denoted as  $Acc_{err}$ . Unlike the previous metric,  $\Delta Prob$  may not affect PLMs’ outputs, despite its suitability for evaluating DKNs’ effect across all queries. Therefore, to measure enhanced DKNs’ impact on PLM performance, we choose  $Acc_{err}$ , as it requires a significant change in  $Prob$  to alter the output. Similarly, we adopt two baselines: enhancing knowledge neurons and enhancing random neurons. Notably, enhancing no neurons directly corresponds to  $Q_{err}^*$ , with an accuracy of 0. Figure 6 presents our results.

**Findings** (1) As depicted in Figure 5, compared to the baselines, suppressing DKNs leads to a significant decrease in PLMs’ prediction probabilities for  $Q^*$ , indicating that suppressing DKNs impairs the robustness of PLMs. (2) Figure 6 shows that, compared to the baselines, enhancing DKNs enables PLMs to provide correct answers to queries that were previously answered incorrectly due to interference, demonstrating that enhancing DKNs improves the robustness of PLMs. In summary, we can conclude that **DKNs can help PLMs cope with input interference**, indicating their contribution to robustness.

## 4.2 Fact-Checking

**Explanation** Since factual errors can also be seen as a form of perturbation, we naturally think of conducting fact-checking experiments based on DKNs.

Method	GPT-2 (Relation)			LLaMA2 (Relation)			GPT-2 (Golden)			LLaMA2 (Golden)		
	P	R	F1	P	R	F1	P	R	F1	P	R	F1
KNs	0.489	0.511	0.500	0.481	0.455	0.468	0.600	0.636	0.618	0.568	0.741	0.643
PLMs	0.015	0.015	0.015	0.035	0.036	0.036	0.015	0.015	0.015	0.035	0.036	0.036
<b>DKNs</b>	<b>0.497</b>	<b>0.520</b>	<b>0.508</b>	<b>0.505</b>	<b>0.500</b>	<b>0.502</b>	<b>0.549</b>	<b>0.848</b>	<b>0.667</b>	<b>0.530</b>	<b>0.955</b>	<b>0.682</b>

Table 2: Results of the fact-checking experiment. ‘‘Relation’’ represents DKNs or KNs based on relation, and ‘‘Golden’’ represents DKNs or KNs obtained using true answers. ‘‘Golden’’ better reflects DKNs’ actual capability, as ‘‘Relation’’ may yield mismatches between DKNs (or KNs) and specific queries due to aggregation and filtering.

Unlike knowledge localization where the true value  $y^*$  is known, fact-checking does not allow us to know  $y^*$  in advance. To address this, we divide the factual knowledge by relation. For queries  $Q^r$  corresponding to a relation, we split them into two parts:  $Q^{r1}$  for DKNs acquisition and  $Q^{r2}$  for testing. For a given query  $Q_i^{r1}$ , the corresponding set of DKNs is denoted as  $D_i^{r1}$ . We aggregate them and select those that appear more than  $\tau_3$  times, denoting as the relation-based DKNs  $D^r$ :

$$D^r = \left\{ n_i \mid n_i \in \bigcup_{D_i^{r1}} \text{and} \sum_{j=1}^{|\bigcup_{D_i^{r1}}|} \mathbb{1}\{n_i = n_j\} > \tau_3 \right\} \quad (7)$$

where  $\sum_{j=1}^{|\bigcup_{D_i^{r1}}|} \mathbb{1}\{n_i = n_j\}$  is the number of occurrences of  $n_i$ . Then, for a query  $Q_i^{r2}$ , we determine the factual correctness by computing the average attribution score of all neurons in  $D^r$ :

$$FC(Q^{r2}) = \begin{cases} \text{True} & \text{if } \sum_{n \in D^r} \text{Score}(n) / |D^r| > \tau_4, \\ \text{False} & \text{otherwise.} \end{cases} \quad (8)$$

where  $\text{Score}(n)$  denotes the activation score of each neuron  $n$  in  $D^r$  (see Appendix B for details),  $\tau_4$  is a threshold.

**Experimental settings** First, we construct a dataset  $Q_f^{r2}$ . For the queries in  $Q_i^{r2}$ , we replace their correct answers with other answers from the same relation (i.e., incorrect answer), and then perform fact-checking on  $Q_f^{r2}$ . Then, we employ two baseline methods: (1) Fact-checking based on KNs, and (2) Direct fact-checking with PLMs providing True or False answers. Besides using relation-based  $D^r$ , we also utilize  $D_i^r$  obtained from true values (‘‘Golden’’). For evaluation metrics, we choose Precision (P), Recall (R), and F1-Score. Table 2 presents the overall results.

**Findings** (1) DKNs have the strongest fact-checking ability. Comparing the P, R, and F1 values of KNs and PLMs, it can be confirmed that DKNs achieve better results under various settings. (2) The ‘‘Golden’’ results in Table 2 outperform

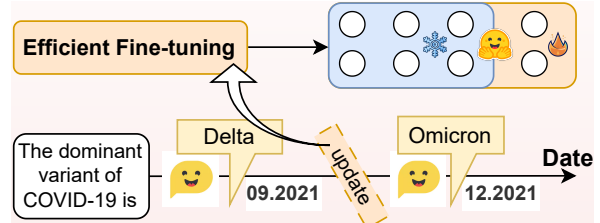


Figure 8: Efficient fine-tuning of PLMs using DKNs to update knowledge across different timestamps.

the results corresponding to ‘‘Relation’’. Thus, the more precise the location of DKNs, the better the fact-checking performance. (3) The fact-checking ability of PLMs is weak. PLMs can correctly answer the query yet fail to identify the fact’s errors due to insufficient fact familiarity. Since true values are used in obtaining  $D^r$ , the fact-checking ability for such facts is stronger.

## 5 The Impact of DKNs on Evolvability

In real-world scenarios, it is meaningful for PLMs to continuously learn or update new knowledge without forgetting old knowledge. Therefore, we investigate one aspect of PLMs’ evolvability, namely the ability of PLMs to learn new knowledge. As shown in Figure 8, given a timestamped query, ‘‘In date \_\_, the dominant variant of COVID-19 is \_\_’’, a PLM may update its answer from ‘‘Delta’’ to ‘‘Omicron’’ across different timestamps, while still retaining the knowledge of the ‘‘Delta’’ variant.

### 5.1 Overlap of DKNs and Parameter Changes

**Experimental settings** We first directly fine-tuning the PLMs and then record the positions of neurons where significant parameter changes occur, denoted as  $\Delta N$ :

$$\Delta N = \{n \mid \Delta P(n) > \tau_{\Delta N}\} \quad (9)$$

$$\Delta P(n) = \sqrt{\left( \frac{\|w_2^{\text{fc}}(n) - w_1^{\text{fc}}(n)\|}{\|w_1^{\text{fc}}(n)\|} \right)^2 + \left( \frac{\|w_2^{\text{proj}}(n) - w_1^{\text{proj}}(n)\|}{\|w_1^{\text{proj}}(n)\|} \right)^2} \quad (10)$$

where  $\tau_{\Delta N}$  is a dynamic threshold,  $\Delta P(n)$  indicates parameter change.  $w_1^{\text{fc}}(n)$  and  $w_1^{\text{proj}}(n)$  sig-

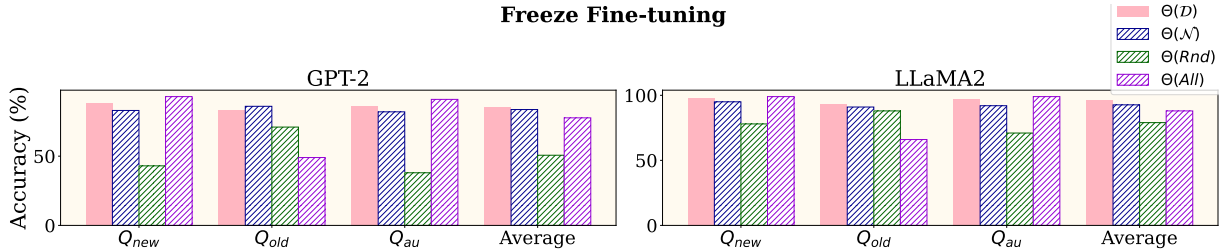


Figure 9: Results of the fine-tuning experiment.  $\Theta$  denotes the areas unfrozen in PLMs during fine-tuning, while  $Q_{new}$ ,  $Q_{old}$ , and  $Q_{au}$  symbolize new, old, and augmented queries, respectively. A method is considered valuable for practical applications only if it achieves high and balanced accuracy across all datasets.

nify the feed-forward and projection weights of neuron  $n$  before fine-tuning, respectively, while  $w_2^{fc}(n)$  and  $w_2^{proj}(n)$  are their post-fine-tuning counterparts. Then, we identify the corresponding  $\mathcal{D}$  through Algorithm 1, and calculate the overlap between  $\mathcal{D}$  and  $\Delta N$ :

$$O(\mathcal{D}, \Delta N) = |\mathcal{D} \cap \Delta N| / |\mathcal{D}| \quad (11)$$

For comparison, we choose the KNs ( $\mathcal{N}$ ) and randomly chosen neurons ( $Rnd$ ) as baselines.

**Findings** The PLMs indeed utilize DKNs to learn new knowledge. Figure 10 reveals  $O(\mathcal{D}, \Delta N)$  is highest, indicating that the parameter change area has the largest overlap degree with DKNs. Notably, LLaMA2 has more irrelevant neurons due to its larger number of parameters, but this does not conflict with its high  $O(\mathcal{D}, \Delta N)$ .

## 5.2 DKNs Guide PLMs to Learn Knowledge

**Experimental settings** Since PLMs primarily utilize DKNs during fine-tuning, we naturally consider leveraging them for efficient fine-tuning to update knowledge (Figure 8). Given a dataset and corresponding  $\mathcal{D}$ , we freeze the parameters outside of  $\mathcal{D}$  during fine-tuning. Then we evaluate the role of DKNs in learning new knowledge by conducting the experiments with three distinct datasets. (1)  $Q_{new}$ : Comprising queries introducing new knowledge, aimed at evaluating PLMs’ grasp of it. (2)  $Q_{old}$ : Comprising knowledge previously mastered by PLMs, with no overlap with  $Q_{new}$ , used to assess if PLMs have forgotten old knowledge. (3)  $Q_{au}$ : Comprising augmented data, which rephrases  $Q_{new}$  into semantically identical but differently expressed queries. We use it to prove that PLMs actually learn knowledge, not just superficial semantic association. Besides unfreezing DKNs ( $\Theta(\mathcal{D})$ ), we select three baseline methods: unfreezing KNs ( $\Theta(\mathcal{N})$ ), unfreezing a random set of neurons equal in number to the DKNs ( $\Theta(Rnd)$ ), and unfreezing all neurons, i.e., direct fine-tuning ( $\Theta(All)$ ).

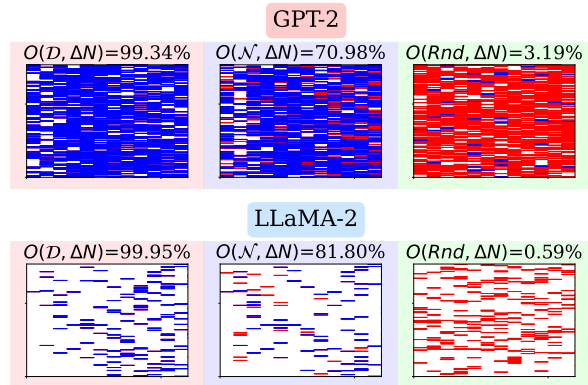


Figure 10: Results of the parameter changes experiment. Blue represents the overlap between neurons and the area of parameter changes, while red indicates the non-overlapping portions. “ $O$ ” is the degree of overlap, with higher values suggest that PLMs rely more on these neurons for learning new knowledge.

**Findings** (1) **DKNs enable PLMs to learn new knowledge more efficiently while preserving old knowledge.** Figure 9 shows that while direct fine-tuning can yield high accuracy on  $Q_{new}$ , it may cause a markedly reduced accuracy on  $Q_{old}$ . However, unfreezing DKNs not only maintains high accuracy on  $Q_{new}$  but also significantly improves performance on  $Q_{old}$  compared to direct fine-tuning. This approach offers a new potential solution to the problem of catastrophic forgetting and makes the fine-tuning process more efficient. (2) Comparing the accuracy of different methods, we find that DKNs possess the strongest ability to guide PLMs in learning new knowledge compared to other neurons. (3) The close accuracy rates between  $Q_{au}$  and  $Q_{new}$  indicate that PLMs can learn genuine knowledge rather than just superficial semantic information.

## 6 The Impact of DKNs on Complexity

### 6.1 Explanation of Complexity

**Hypotheses** In auto-regressive models, complexity, denoted by  $\mathcal{C}(M)$ , is determined by the num-

ber of model parameters (Hu et al., 2021). Given two PLMs,  $M_1$  and  $M_2$ . Let  $\mathcal{D}^*(M)$  be the PLMs’ degeneracy,  $\mathcal{T}$  be a downstream task using DKNs. The performance of  $M$  on  $\mathcal{T}$  is indicated by  $\mathcal{P}(M, \mathcal{T})$ , and  $Param(M)$  is the number of parameters in  $M$ . We propose two hypotheses:

**Hypothesis 1**  $Param(M_1) < Param(M_2) \Rightarrow \mathcal{C}(M_1) < \mathcal{C}(M_2)$

**Hypothesis 2**  $\mathcal{P}(M_1, \mathcal{T}) < \mathcal{P}(M_2, \mathcal{T}) \Rightarrow \mathcal{D}^*(M_1) < \mathcal{D}^*(M_2)$  (12)

**Findings** (1) When DKNs are suppressed, LLaMA2 experiences a more significant drop in prediction probability compared to GPT-2. Conversely, enhancing DKNs leads to a more notable increase in LLaMA2’s accuracy, as illustrated in Figures 5 and 6. (2) LLaMA2 surpasses GPT-2 in fact-checking abilities when utilizing DKNs, as indicated by Table 2. (3) Compared to GPT-2, LLaMA2 has a stronger ability to learn new knowledge, retain old knowledge, and acquire genuine knowledge, as detailed in Figure 9. The aforementioned phenomena indicate that  $\mathcal{P}(\text{GPT-2}, \mathcal{T}) < \mathcal{P}(\text{LLaMA2}, \mathcal{T})$ , and given that  $Param(\text{GPT-2}) < Param(\text{LLaMA2})$ , according to Hypotheses 12, we can conclude:

$$\mathcal{C}(\text{GPT-2}) < \mathcal{C}(\text{LLaMA2}) \Rightarrow \mathcal{D}^*(\text{GPT-2}) < \mathcal{D}^*(\text{LLaMA2}) \quad (13)$$

In summary, **degeneracy is positively correlated with complexity**.

## 6.2 Fact-Checking within Complex Texts

**Experimental settings** We investigate the potential of DKNs when integrated with PLMs’ other abilities, through an experiment that melds the model’s ability to comprehend complex texts with DKNs’ fact-checking ability within such texts. Given a triplet-form query  $Q$ , first, we rewrite it into a complex text  $Q^c$ . To perform fact-checking on  $Q^c$ , the PLMs must possess both reasoning ability and fact-checking ability. Then, we provide a prompt to PLMs, requiring them to first identify both the query and answer in  $Q^c$ , and then perform fact-checking according to Equation 8.

**Findings** (1) PLMs with the ability to understand complex texts can employ DKNs for fact-checking within complex texts. Table 3 shows LLaMA2’s ability to utilize DKNs for fact-checking. However, GPT-2 only echoes prompts without yielding meaningful results, thus its exclusion from the table. (2) When DKNs’ locations are imprecise, pre-understanding complex texts can enhance fact-checking ability. The results of “Relation” in Table

Method	LLaMA2 (Relation)			LLaMA2 (Golden)		
	P	R	F1	P	R	F1
KNs	0.59	0.48	0.53	0.64	0.43	0.51
PLMs	0.02	0.02	0.02	0.02	0.02	0.02
<b>DKNs</b>	<b>0.54</b>	<b>0.58</b>	<b>0.56</b>	<b>0.52</b>	<b>0.78</b>	<b>0.62</b>

Table 3: Results of the fact-checking in complex texts.

3 surpass those in Table 2 (F1: 0.56 to 0.502), but the results of “Golden” decrease (F1: 0.62 to 0.682). This indicates that relation-based DKNs lack precision, but complex text understanding increases PLMs’ sensitivity to specific facts.

## 7 Related Work

Petroni et al. (2019b) argue that numerous factual knowledge exists within PLMs and suggest using “fill-in-the-blank” cloze tasks determine if the models have grasped specific facts. Meanwhile, Geva et al. (2021) suggest that MLP neurons within transformer models operate as key-value memories. Building on this, Dai et al. (2022) uncover that some MLP neurons are capable of storing factual knowledge, termed as knowledge neurons (KNs). Lundstrom et al. (2022) confirm the reliability of their knowledge localization method, while subsequent knowledge editing experiments by Meng et al. (2022) and Meng et al. (2023) reinforce that MLP neurons indeed store factual knowledge. Moreover, Geva et al. (2023) investigate KN dynamics, and Chen et al. (2024a) discover multiple distinct KN sets storing identical facts, termed degenerate knowledge neurons (DKNs). Additionally, researchers have discovered various other KNs with unique properties (Wang et al., 2022; Tang et al., 2024). In summary, despite limitations (Niu et al., 2024; Bricken et al., 2023; Chen et al., 2024b), the KN-based analysis approach remains valuable for further research.

## 8 Conclusion

This paper delves into the mechanisms of factual knowledge storage in PLMs. First, we provide a comprehensive definition of DKNs that covers both structural and functional aspects, pioneering the study of the internal structures of PLMs. Based on this, we introduce the neurological topology clustering method for more precise DKN acquisition. Finally, through the lens of DKNs, we investigate the robustness, evolvability, and complexity of PLMs. Extensive experiments demonstrate the critical role of DKNs in PLMs.



## 547 **Limitations**

548 First, limited by computational resources, our study  
549 involves only the GPT-2 and LLaMA2-7b models.  
550 To further validate the scalability of our methods  
551 and conclusions, it is imperative to conduct studies  
552 on larger models. Second, our research is confined  
553 to factual knowledge, and whether similar findings  
554 apply to other types of knowledge remains to be  
555 investigated. Finally, the generalizability of our  
556 findings across different languages and cultural con-  
557 texts remains an open question. Our study utilizes  
558 datasets in English, limiting our ability to assess  
559 the performance and applicability of our methods  
560 across non-English languages and datasets that en-  
561 compass a broader range of cultural knowledge.

## 562 **Ethics Statement**

563 Our research aims to deepen the understanding  
564 and enhance the functionality of PLMs by inves-  
565 tigating the role of DKNs. While our research  
566 seeks to enhance the understanding and function-  
567 ality of PLMs by exploring the role of DKNs, we  
568 are acutely aware of the potential for misuse of  
569 these findings. The increased capabilities of PLMs,  
570 driven by insights into DKNs, could potentially be  
571 exploited for generating misleading information,  
572 manipulating public opinion, or other malicious  
573 purposes. It is not our intention to facilitate such  
574 activities. Instead, our goal is to contribute to the  
575 scientific community’s knowledge base, enabling  
576 the development of more robust, accurate, and ethi-  
577 cally aligned language technologies..

578 We emphasize the responsibility of using these  
579 insights ethically, to avoid the exploitation of en-  
580 hanced PLM capabilities for harmful purposes. To  
581 this end, we advocate for transparency in research,  
582 collaboration across disciplines for ethical over-  
583 sight, and the development of regulatory frame-  
584 works to ensure that advancements in PLM tech-  
585 nology contribute positively to society.

## 586 **References**

587 Mohiuddin Ahmed, Raihan Seraj, and Syed Mo-  
588 hammed Shamsul Islam. 2020. **The k-means algo-**  
589 **rithm: A comprehensive survey and performance**  
590 **evaluation.** *Electronics*, 9(8):1295.

591 Zeyuan Allen-Zhu and Yuanzhi Li. 2023. **Physics of**  
592 **language models: Part 3.2, knowledge manipulation.**  
593 *ArXiv preprint*, abs/2309.14402.

Serguei Barannikov. 1994. **The framed morse complex**  
594 **and its invariants.** *Advances in Soviet Mathematics*,  
595 21:93–116. 596

Trenton Bricken, Adly Templeton, Joshua Batson,  
597 Brian Chen, Adam Jermy, Tom Conerly, Nick  
598 Turner, Cem Anil, Carson Denison, Amanda Askill,  
599 Robert Lasenby, Yifan Wu, Shauna Kravec, Nicholas  
600 Schiefer, Tim Maxwell, Nicholas Joseph, Zac  
601 Hatfield-Dodds, Alex Tamkin, Karina Nguyen,  
602 Brayden McLean, Josiah E Burke, Tristan Hume,  
603 Shan Carter, Tom Henighan, and Christopher  
604 Olah. 2023. **Towards monosemanticity: Decom-**  
605 **posing language models with dictionary learning.**  
606 *Transformer Circuits Thread.* [https://transformer-](https://transformer-circuits.pub/2023/monosemantic-features/index.html)  
607 [circuits.pub/2023/monosemantic-](https://transformer-circuits.pub/2023/monosemantic-features/index.html)  
608 [features/index.html.](https://transformer-circuits.pub/2023/monosemantic-features/index.html) 609

Gunnar Carlsson. 2009. **Topology and data.** *Bulletin of*  
610 *the American Mathematical Society*, 46(2):255–308. 611

Tianyi Chen, Yeliz Yesilada, and Simon Harper. 2010.  
612 **What input errors do you experience? typing and**  
613 **pointing errors of mobile web users.** *International*  
614 *journal of human-computer studies*, 68(3):138–157. 615

Yuheng Chen, Pengfei Cao, Yubo Chen, Kang Liu, and  
616 Jun Zhao. 2024a. **Journey to the center of the knowl-**  
617 **edge neurons: Discoveries of language-independent**  
618 **knowledge neurons and degenerate knowledge neu-**  
619 **rons.** In *Proceedings of the AAAI Conference on*  
620 *Artificial Intelligence.* 621

Yuheng Chen, Pengfei Cao, Yubo Chen, Kang Liu, and  
622 Jun Zhao. 2024b. **Knowledge localization: Mis-**  
623 **sion not accomplished? enter query localization!**  
624 *Preprint*, arXiv:2405.14117. 625

David Cohen-Steiner, Herbert Edelsbrunner, and John  
626 Harer. 2005. **Stability of persistence diagrams.** In  
627 *Proceedings of the twenty-first annual symposium on*  
628 *Computational geometry*, pages 263–271. 629

Damai Dai, Li Dong, Yaru Hao, Zhifang Sui, Baobao  
630 Chang, and Furu Wei. 2022. **Knowledge neurons in**  
631 **pretrained transformers.** In *Proceedings of the 60th*  
632 *Annual Meeting of the Association for Computational*  
633 *Linguistics (Volume 1: Long Papers)*, pages 8493–  
634 8502, Dublin, Ireland. Association for Computational  
635 Linguistics. 636

Tamal K Dey, Dayu Shi, and Yusu Wang. 2019. **Simba:**  
637 **An efficient tool for approximating rips-filtration per-**  
638 **sistence via simplicial batch collapse.** *Journal of*  
639 *Experimental Algorithmics (JEA)*, 24:1–16. 640

Bhuwan Dhingra, Jeremy R. Cole, Julian Martin  
641 Eisenschlos, Daniel Gillick, Jacob Eisenstein, and  
642 William W. Cohen. 2022. **Time-aware language mod-**  
643 **els as temporal knowledge bases.** *Transactions of the*  
644 *Association for Computational Linguistics*, 10:257–  
645 273. 646

Gerald M Edelman and Joseph A Gally. 2001. **Degen-**  
647 **eracy and complexity in biological systems.** *Pro-*  
648 *ceedings of the National Academy of Sciences*,  
649 98(24):13763–13768. 650

651	Herbert Edelsbrunner, John Harer, et al. 2008. Persistent homology—a survey. <i>Contemporary mathematics</i> , 453(26):257–282.	704
652		705
653		706
654	Herbert Edelsbrunner and John L Harer. 2022. <i>Computational topology: an introduction</i> . American Mathematical Society.	707
655		708
656		709
657	Martin Ester, Hans-Peter Kriegel, Jörg Sander, Xiaowei Xu, et al. 1996. A density-based algorithm for discovering clusters in large spatial databases with noise. In <i>kdd</i> , volume 96, pages 226–231.	710
658		711
659		712
660		713
661	Jean-Claude Fernandez, Laurent Mounier, and Cyril Pachon. 2005. A model-based approach for robustness testing. In <i>Testing of Communicating Systems: 17th IFIP TC6/WG 6.1 International Conference, TestCom 2005, Montreal, Canada, May 31-June, 2005. Proceedings 17</i> , pages 333–348. Springer.	714
662		715
663		716
664		717
665		718
666		
667	Mor Geva, Jasmijn Bastings, Katja Filippova, and Amir Globerson. 2023. Dissecting recall of factual associations in auto-regressive language models. <i>Preprint</i> , arXiv:2304.14767.	719
668		720
669		721
670		722
671	Mor Geva, Roei Schuster, Jonathan Berant, and Omer Levy. 2021. Transformer feed-forward layers are key-value memories. In <i>Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing</i> , pages 5484–5495, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.	723
672		724
673		725
674		726
675		
676		
677		
678	Zhijiang Guo, Michael Schlichtkrull, and Andreas Vlachos. 2022. A survey on automated fact-checking. <i>Transactions of the Association for Computational Linguistics</i> , 10:178–206.	727
679		728
680		729
681		730
682	Xia Hu, Lingyang Chu, Jian Pei, Weiqing Liu, and Jiang Bian. 2021. Model complexity of deep learning: A survey. <i>Preprint</i> , arXiv:2103.05127.	731
683		732
684		733
685	Nora Kassner, Philipp Dufter, and Hinrich Schütze. 2021. Multilingual LAMA: Investigating knowledge in multilingual pretrained language models. In <i>Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume</i> , pages 3250–3258, Online. Association for Computational Linguistics.	734
686		735
687		736
688		737
689		738
690		739
691		740
692	Michael Kerber and Raghendra Sharathkumar. 2013. Approximate čech complex in low and high dimensions. In <i>International Symposium on Algorithms and Computation</i> , pages 666–676. Springer.	741
693		742
694		743
695		744
696	Marc Kirschner and John Gerhart. 1998. Evolvability. <i>Proceedings of the National Academy of Sciences</i> , 95(15):8420–8427.	745
697		746
698		747
699	Daniel D Lundstrom, Tianjian Huang, and Meisam Razaviyayn. 2022. A rigorous study of integrated gradients method and extensions to internal neuron attributions. In <i>International Conference on Machine Learning</i> , pages 14485–14508. PMLR.	748
700		749
701		750
702		751
703		752
		753
		754
		755
		756
		757
		758
		759
		760
		761
	Mourad Mars. 2022. From word embeddings to pre-trained language models: A state-of-the-art walk-through. <i>Applied Sciences</i> , 12(17):8805.	
	Paul H Mason. 2015. Degeneracy: Demystifying and destigmatizing a core concept in systems biology. <i>Complexity</i> , 20(3):12–21.	
	Kevin Meng, David Bau, Alex J Andonian, and Yonatan Belinkov. 2022. Locating and editing factual associations in GPT. In <i>Advances in Neural Information Processing Systems</i> .	
	Kevin Meng, Arnab Sen Sharma, Alex J Andonian, Yonatan Belinkov, and David Bau. 2023. Mass-editing memory in a transformer. In <i>The Eleventh International Conference on Learning Representations</i> .	
	Fionn Murtagh and Pedro Contreras. 2012. Algorithms for hierarchical clustering: an overview. <i>Wiley Interdisciplinary Reviews: Data Mining and Knowledge Discovery</i> , 2(1):86–97.	
	Jingcheng Niu, Andrew Liu, Zining Zhu, and Gerald Penn. 2024. What does the knowledge neuron thesis have to do with knowledge? In <i>The Twelfth International Conference on Learning Representations</i> .	
	OpenAI, :, Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altschmidt, Sam Altman, Shyamal Anadkat, Red Avila, Igor Babuschkin, Suchir Balaji, Valerie Balcom, Paul Baltescu, Haiming Bao, Mo Bavarian, Jeff Belgum, Irwan Bello, Jake Berdine, Gabriel Bernadett-Shapiro, Christopher Berner, Lenny Bogdonoff, Oleg Boiko, Madelaine Boyd, Anna-Luisa Brakman, Greg Brockman, Tim Brooks, Miles Brundage, Kevin Button, Trevor Cai, Rosie Campbell, Andrew Cann, Brittany Carey, Chelsea Carlson, Rory Carmichael, Brooke Chan, Che Chang, Fotis Chantzis, Derek Chen, Sully Chen, Ruby Chen, Jason Chen, Mark Chen, Ben Chess, Chester Cho, Casey Chu, Hyung Won Chung, Dave Cummings, Jeremiah Currier, Yunxing Dai, Cory Decareaux, Thomas Degry, Noah Deutsch, Damien Deville, Arka Dhar, David Dohan, Steve Dowling, Sheila Dunning, Adrien Ecoffet, Atty Eleti, Tyna Eloundou, David Farhi, Liam Fedus, Niko Felix, Simón Posada Fishman, Juston Forte, Isabella Fulford, Leo Gao, Elie Georges, Christian Gibson, Vik Goel, Tarun Gogineni, Gabriel Goh, Rapha Gontijo-Lopes, Jonathan Gordon, Morgan Grafstein, Scott Gray, Ryan Greene, Joshua Gross, Shixiang Shane Gu, Yufei Guo, Chris Hallacy, Jesse Han, Jeff Harris, Yuchen He, Mike Heaton, Johannes Heidecke, Chris Hesse, Alan Hickey, Wade Hickey, Peter Hoeschele, Brandon Houghton, Kenny Hsu, Shengli Hu, Xin Hu, Joost Huizinga, Shantanu Jain, Shawn Jain, Joanne Jang, Angela Jiang, Roger Jiang, Haozhun Jin, Denny Jin, Shino Jomoto, Billie Jonn, Heewoo Jun, Tomer Kaftan, Łukasz Kaiser, Ali Kamali, Ingmar Kanitscheider, Nitish Shirish Keskar, Tabarak Khan, Logan Kilpatrick, Jong Wook	

762	Kim, Christina Kim, Yongjik Kim, Hendrik Kirchner, Jamie Kiros, Matt Knight, Daniel Kokotajlo, Lukasz Kondraciuk, Andrew Kondrich, Aris Konstantinidis, Kyle Kosic, Gretchen Krueger, Vishal Kuo, Michael Lampe, Ikai Lan, Teddy Lee, Jan Leike, Jade Leung, Daniel Levy, Chak Ming Li, Rachel Lim, Molly Lin, Stephanie Lin, Mateusz Litwin, Theresa Lopez, Ryan Lowe, Patricia Lue, Anna Makanju, Kim Malfacini, Sam Manning, Todor Markov, Yaniv Markovski, Bianca Martin, Katie Mayer, Andrew Mayne, Bob McGrew, Scott Mayer McKinney, Christine McLeavey, Paul McMillan, Jake McNeil, David Medina, Aalok Mehta, Jacob Menick, Luke Metz, Andrey Mishchenko, Pamela Mishkin, Vinnie Monaco, Evan Morikawa, Daniel Mossing, Tong Mu, Mira Murati, Oleg Murk, David Mély, Ashvin Nair, Reiichiro Nakano, Rajeesh Nayak, Arvind Neelakantan, Richard Ngo, Hyeonwoo Noh, Long Ouyang, Cullen O’Keefe, Jakub Pachocki, Alex Paino, Joe Palermo, Ashley Pantuliano, Giambattista Parascandolo, Joel Parish, Emy Parparita, Alex Passos, Mikhail Pavlov, Andrew Peng, Adam Perelman, Filipe de Avila Belbute Peres, Michael Petrov, Henrique Ponde de Oliveira Pinto, Michael, Pokorny, Michelle Pokrass, Vitchyr Pong, Tolly Powell, Alethea Power, Boris Power, Elizabeth Proehl, Raul Puri, Alec Radford, Jack Rae, Aditya Ramesh, Cameron Raymond, Francis Real, Kendra Rimbach, Carl Ross, Bob Rotsted, Henri Roussez, Nick Ryder, Mario Saltarelli, Ted Sanders, Shibani Santurkar, Girish Sastry, Heather Schmidt, David Schnurr, John Schulman, Daniel Selsam, Kyla Sheppard, Toki Sherbakov, Jessica Shieh, Sarah Shoker, Pranav Shyam, Szymon Sidor, Eric Sigler, Maddie Simens, Jordan Sitkin, Katarina Slama, Ian Sohl, Benjamin Sokolowsky, Yang Song, Natalie Staudacher, Felipe Petroski Such, Natalie Summers, Ilya Sutskever, Jie Tang, Nikolas Tezak, Madeleine Thompson, Phil Tillet, Amin Tootoonchian, Elizabeth Tseng, Preston Tuggle, Nick Turley, Jerry Tworek, Juan Felipe Cerón Uribe, Andrea Vallone, Arun Vijayvergiya, Chelsea Voss, Carroll Wainwright, Justin Jay Wang, Alvin Wang, Ben Wang, Jonathan Ward, Jason Wei, CJ Weinmann, Akila Welihinda, Peter Welinder, Jiayi Weng, Lilian Weng, Matt Wiethoff, Dave Willner, Clemens Winter, Samuel Wolrich, Hannah Wong, Lauren Workman, Sherwin Wu, Jeff Wu, Michael Wu, Kai Xiao, Tao Xu, Sarah Yoo, Kevin Yu, Qiming Yuan, Wojciech Zaremba, Rowan Zellers, Chong Zhang, Marvin Zhang, Shengjia Zhao, Tianhao Zheng, Juntang Zhuang, William Zhuk, and Barret Zoph. 2023. <a href="#">Gpt-4 technical report</a> . <i>Preprint</i> , arXiv:2303.08774.	
763		
764		
765		
766		
767		
768		
769		
770		
771		
772		
773		
774		
775		
776		
777		
778		
779		
780		
781		
782		
783		
784		
785		
786		
787		
788		
789		
790		
791		
792		
793		
794		
795		
796		
797		
798		
799		
800		
801		
802		
803		
804		
805		
806		
807		
808		
809		
810		
811		
812		
813		
814		
815	Nina Otter, Mason A Porter, Ulrike Tillmann, Peter Grindrod, and Heather A Harrington. 2017. <a href="#">A roadmap for the computation of persistent homology</a> . <i>EPJ Data Science</i> , 6:1–38.	
816		
817		
818		
819	Fabio Petroni, Patrick Lewis, Aleksandra Piktus, Tim Rocktäschel, Yuxiang Wu, Alexander H. Miller, and Sebastian Riedel. 2020. <a href="#">How context affects language models’ factual predictions</a> . In <i>Automated Knowledge Base Construction</i> .	
820		
821		
822		
823		
	Fabio Petroni, Tim Rocktäschel, Sebastian Riedel, Patrick Lewis, Anton Bakhtin, Yuxiang Wu, and Alexander Miller. 2019a. <a href="#">Language models as knowledge bases?</a> In <i>Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)</i> , pages 2463–2473, Hong Kong, China. Association for Computational Linguistics.	824 825 826 827 828 829 830 831 832
	Fabio Petroni, Tim Rocktäschel, Sebastian Riedel, Patrick Lewis, Anton Bakhtin, Yuxiang Wu, and Alexander Miller. 2019b. <a href="#">Language models as knowledge bases?</a> In <i>Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)</i> , pages 2463–2473, Hong Kong, China. Association for Computational Linguistics.	833 834 835 836 837 838 839 840 841
	Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. 2019. <a href="#">Language models are unsupervised multitask learners</a> . <i>OpenAI blog</i> , 1(8):9.	842 843 844 845
	Tianyi Tang, Wenyang Luo, Haoyang Huang, Dongdong Zhang, Xiaolei Wang, Xin Zhao, Furu Wei, and Ji-Rong Wen. 2024. <a href="#">Language-specific neurons: The key to multilingual capabilities in large language models</a> . <i>Preprint</i> , arXiv:2402.16438.	846 847 848 849 850
	Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shrutit Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony Hartshorn, Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel Kloumann, Artem Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov, Pushkar Mishra, Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan Saladi, Alan Schelten, Ruan Silva, Eric Michael Smith, Ranjan Subramanian, Xiaoqing Ellen Tan, Binh Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zheng Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurelien Rodriguez, Robert Stojnic, Sergey Edunov, and Thomas Scialom. 2023. <a href="#">Llama 2: Open foundation and fine-tuned chat models</a> . <i>Preprint</i> , arXiv:2307.09288.	851 852 853 854 855 856 857 858 859 860 861 862 863 864 865 866 867 868 869 870 871 872 873
	Alessandro Verri, Claudio Uras, Patrizio Frosini, and Massimo Ferri. 1993. <a href="#">On the use of size functions for shape analysis</a> . <i>Biological cybernetics</i> , 70(2):99–107.	874 875 876 877
	Xiaozhi Wang, Kaiyue Wen, Zhengyan Zhang, Lei Hou, Zhiyuan Liu, and Juanzi Li. 2022. <a href="#">Finding skill neurons in pre-trained transformer-based language models</a> . In <i>Conference on Empirical Methods in Natural Language Processing</i> .	878 879 880 881 882

James Whitacre and Axel Bender. 2010. [Degeneracy: a design principle for achieving robustness and evolvability](#). *Journal of theoretical biology*, 263(1):143–153.

James M Whitacre. 2010. [Degeneracy: a link between evolvability, robustness and complexity in biological systems](#). *Theoretical Biology and Medical Modelling*, 7:1–17.

Zeyuan Allen Zhu and Yuanzhi Li. 2023. [Physics of language models: Part 3.1, knowledge storage and extraction](#). *ArXiv preprint*, abs/2309.14316.

Afra Zomorodian and Gunnar Carlsson. 2004. [Computing persistent homology](#). In *Proceedings of the twentieth annual symposium on Computational geometry*, pages 347–356.

## A Experimental Hyperparameters

### A.1 Hardware specification and environment.

We ran our experiments on the machine equipped with the following specifications:

- CPU: Intel(R) Xeon(R) CPU E5-2680 v4 @ 2.40GHz, Total CPUs: 56
- GPU: NVIDIA GeForce RTX 3090, 24576 MiB (10 units)
- Software:
  - Python Version: 3.10.10
  - PyTorch Version: 2.0.0+cu117

### A.2 Experimental Hyperparameters of Neurological Topology Clustering

In Section 3, where we obtain degenerate knowledge neurons, the primary hyperparameters are  $\tau_1$  and  $\tau_2$ . First,  $\tau_1$  is a dynamic threshold, set as

$$\tau_1 = 0.5 \times \max(R_{\text{persist}}(\mathcal{B}_1), \dots, R_{\text{persist}}(\mathcal{B}_n)) \quad (14)$$

Then,  $\tau_2$  is a fixed value.

$$\tau_2 = 0.3 \quad (15)$$

In this experiment, some data led to excessively large changes in predictive probability, indicating that the PLMs had not originally mastered this factual knowledge, resulting in a very low initial predictive probability. To study the storage mechanism of factual knowledge, it’s essential to investigate facts already grasped by the model. Therefore, we set a threshold to exclude data that caused extreme

changes in predictive probability. If the change in predictive probability  $\Delta P_{\text{Prob}}$  satisfies:

$$\Delta P_{\text{Prob}} > 900 \quad (16)$$

then that data is excluded.

### A.3 Experimental Hyperparameters of Degneracy and Robustness

In Section 4, the first experiment under Query-Perturbation, namely the Suppressing DKN experiment, similar to A.2, excludes data that satisfies the condition:

$$Prob_{\text{sup}} > 900 \quad (17)$$

In the Fact-Checking experiment,  $\tau_3$  is similar to  $\tau_1$  as a dynamic threshold. We first count the total number of neurons, then set

$$\tau_3 = 0.7 \times N_{\text{total}} \quad (18)$$

where  $N_{\text{total}}$  represents the total number of neurons.

### A.4 Experimental Hyperparameters of Degeneracy and Evolvability

In Section 5, the experimental hyperparameter  $\tau_{\Delta N}$  for the Overlap of DKN and Parameter Changes experiment is set as a dynamic threshold. The process involves calculating the maximum value of  $\Delta P(n)$  according to Equation 10. Once this value is determined,  $\tau_{\Delta N}$  is set differently based on the model in use. For the GPT-2 model, the threshold  $\tau_{\Delta N}$  is calculated as:

$$\tau_{\Delta N} = 0.04 \times \max(\Delta P(n_1), \Delta P(n_2), \dots, \Delta P(n_k)) \quad (19)$$

In contrast, for the Llama2 model, the calculation of the threshold  $\tau_{\Delta N}$  is slightly adjusted:

$$\tau_{\Delta N} = 0.05 \times \max(\Delta P(n_1), \Delta P(n_2), \dots, \Delta P(n_k)) \quad (20)$$

This distinction in the calculation of  $\tau_{\Delta N}$  reflects the specific characteristics and performance considerations of each model.

### A.5 Experimental Hyperparameters of Degneracy and Complexity

In Section 6, during our fact-checking experiments on complex texts in Section 6, we provide PLMs with a prompt requiring them to first extract factual knowledge in the form of triples and then perform fact-checking using DKNs. Our prompt is as follows:

Given a complex statement, simplify it into a straightforward sentence that clearly states the subject, predicate, and object. The simplified sentence should capture the key information from the complex statement in a concise format. Here are a few examples to illustrate the conversion:

1. Complex:  $\_X\_$ , renowned for their career in sports, plays for a team or club. This association is a significant part of their professional journey, reflecting their skill and dedication in their field.

Simplified: Valentino Rossi plays for  $\_X\_$ .

2. Complex:  $\_X\_$ , a key figure in their domain, holds an important position. This role is marked by responsibilities and influence, shaping the direction and success of their organization or country.

Simplified: Robert Mugabe holds the position of  $\_X\_$ .

3. Complex:  $\_X\_$ , known for their intellectual and professional achievements, attended a notable institution. Their time at this institution played a crucial role in shaping their career and perspectives.

Simplified: Malala Yousafzai attended  $\_X\_$ .

4. Complex: As the chair of  $\_X\_$ , this individual is a prominent figure in the business world. This position highlights their leadership and strategic vision in guiding the company towards growth and innovation.

Simplified:  $\_X\_$  is the chair of bp.

5. Complex:  $\_X\_$ , a figure of significant influence and ideals, is a member of a notable group or organization. Their membership reflects their commitment to certain values and goals, contributing to the group's impact.

Simplified: Alexei Navalny is a member of the  $\_X\_$ .

6. Complex:  $\_X\_$  works for a notable company or organization, playing a pivotal role in driving innovation and excellence in their field.

Simplified: Elon Musk works for  $\_X\_$ .

7. Complex: As the head of the government of  $\_X\_$ , this political figure plays a central role in shaping the nation's policies and international relations.

Simplified:  $\_X\_$  is the head of the government of France.

8. Complex: The head coach of  $\_X\_$  is known for their expertise in coaching, with leadership and strategic skills crucial in steering the team to success.

Simplified:  $\_X\_$  is the head coach of Philadelphia Eagles.

9. Complex:  $\_X\_$ , a significant player in its sector, is owned by a larger entity. This ownership structure is key to understanding the company's operations and market influence.

Simplified: Google is owned by  $\_X\_$ .

Now, given the following complex statement, apply the same process to simplify it:

Figure 11: The prompt of fact-checking experiments on complex texts.

## A.6 Fact-Checking of Complex Texts in Section 6

During our fact-checking experiments on complex texts in Section 6, we prompt PLMs to simplify complex statements into straightforward sentences. Our prompt is shown in Figure 11.

## B Knowledge Localization

This section introduces the method we use to acquire knowledge neurons. We employ the approach proposed by Chen et al.(2024a), which we will detail below.

Given a query  $q$ , we can define the probability of the correct answer predicted by a PLMs as follows:

$$F(\hat{w}_j^{(l)}) = p(y^* | q, w_j^{(l)} = \hat{w}_j^{(l)}) \quad (21)$$

Here,  $y^*$  represents the correct answer,  $w_j^{(l)}$  denotes the  $j$ -th neuron in the  $l$ -th layer, and  $\hat{w}_j^{(l)}$  is the specific value assigned to  $w_j^{(l)}$ . To calculate the

attribution score for each neuron, we employ the technique of integrated gradients.

To compute the attribution score of a neuron  $w_j^{(l)}$ , we consider the following formulation:

$$\text{Attr}(w_j^{(l)}) = (\bar{w}_j^{(l)} - w_j^{(l)}) \int_0^1 \frac{\partial F(w_j^{(l)} + \alpha(\bar{w}_j^{(l)} - w_j^{(l)}))}{\partial w_j^{(l)}} d\alpha \quad (22)$$

Here,  $\bar{w}_j^{(l)}$  represents the actual value of  $w_j^{(l)}$ ,  $w_j^{(l)}$  serves as the baseline vector for  $w_j^{(l)}$ . The term  $\frac{\partial F(w_j^{(l)} + \alpha(w_j^{(l)} - w_j^{(l)}))}{\partial w_j^{(l)}}$  computes the gradient with respect to  $w_j^{(l)}$ .

Next, we aim to obtain  $w_j^{(l)}$ . Starting from the sentence  $q$ , we acquire a baseline sentence and then encode this sentence as a vector.

Let the baseline sentence corresponding to  $q_i$  be  $q'_i$ , and  $q'_i$  consists of  $m$  words, maintaining a length consistent with  $q$ , denoted as  $q'_i = (q'_{i1} \dots q'_{ik} \dots q'_{im})$ . Since we are using autoregressive models, according to Chen et al.(2024a),

method,  $q'_{ik} = \langle \text{eos} \rangle$ , where  $\langle \text{eos} \rangle$  represents “end of sequence” in auto-regressive models.

The attribution score  $Attr_i(w_j^{(l)})$  for each neuron, given the input  $q_i$ , can be determined using Equation (22). For the computation of the integral, the Riemann approximation method is employed:

$$Attr_i(w_j^l) \approx \frac{\bar{w}_j^{(l)}}{N} \sum_{k=1}^N \frac{\partial F(w_j^{(l)} + \frac{k}{N} \times (\bar{w}_j^{(l)} - w_j^{(l)}))}{\partial w_j^{(l)}} \quad (23)$$

where  $N$  is the number of approximation steps.

Then, the attribution scores for each word  $q_i$  are aggregated and subsequently normalized:

$$Attr(w_j^l) = \frac{\sum_{i=1}^m Attr_i(w_j^l)}{\sum_{j=1}^n \sum_{i=1}^m Attr_i(w_j^l)}, \quad (24)$$

Let  $\mathcal{N}$  be the set of neurons classified as knowledge neurons based on their attribution scores exceeding a predetermined threshold  $\tau$ , for a given input  $q$ . This can be formally defined as:

$$\mathcal{N} = \left\{ w_j^{(l)} \mid Attr(w_j^{(l)}) > \tau \right\} \quad (25)$$

where  $l$  encompassing all layers and  $j$  including all neurons within each layer.

## C Persistent homology

Persistent homology is a method for computing topological features of a space at different spatial resolutions. More persistent features are detected over a wide range of spatial scales and are deemed more likely to represent true features of the underlying space rather than artifacts of sampling, noise, or particular choice of parameters (Carlsson, 2009).

To find the persistent homology of a space, the space must first be represented as a simplicial complex. A distance function on the underlying space corresponds to a filtration of the simplicial complex, that is a nested sequence of increasing subsets. One common method of doing this is via taking the sublevel filtration of the distance to a point cloud, or equivalently, the offset filtration on the point cloud and taking its nerve in order to get the simplicial filtration known as Čech filtration (Kerber and Sharathkumar, 2013). A similar construction uses a nested sequence of Vietoris–Rips complexes known as the Vietoris–Rips filtration (Dey et al., 2019).

### C.1 Definition

In persistent homology, formally, we consider a real-valued function defined on a simplicial complex, denoted as  $f : K \rightarrow \mathbb{R}$ . This function is

required to be non-decreasing on increasing sequences of faces, meaning that for any two faces  $\sigma$  and  $\tau$  in  $K$ , if  $\sigma$  is a face of  $\tau$ , then  $f(\sigma) \leq f(\tau)$ .

For every real number  $a$ , the sublevel set  $K_a = f^{-1}((-\infty, a])$  forms a subcomplex of  $K$ . The values of  $f$  on the simplices in  $K$  create an ordering of these sublevel complexes, which leads to a filtration:

$$\emptyset = K_0 \subseteq K_1 \subseteq \dots \subseteq K_n = K \quad (26)$$

Within this filtration, for  $0 \leq i \leq j \leq n$ , the inclusion  $K_i \hookrightarrow K_j$  induces a homomorphism on the simplicial homology groups for each dimension  $p$ , noted as  $f_p^{i,j} : H_p(K_i) \rightarrow H_p(K_j)$ . The  $p^{\text{th}}$  persistent homology groups are the images of these homomorphisms, and the  $p^{\text{th}}$  persistent Betti numbers  $\beta_p^{i,j}$  are defined as the ranks of these groups (Edelsbrunner and Harer, 2022). Persistent Betti numbers for  $p = 0$  coincide with the size function, an earlier concept related to persistent homology (Verri et al., 1993).

The concept extends further to any filtered complex over a field  $F$ . Such a complex can be transformed into its canonical form, which is a direct sum of filtered complexes of two types: one-dimensional complexes with trivial differential (expressed as  $d(e_{t_i}) = 0$ ) and two-dimensional complexes with trivial homology (expressed as  $d(e_{s_j+r_j}) = e_{r_j}$ ) (Barannikov, 1994).

A persistence module over a partially ordered set  $P$  consists of a collection of vector spaces  $U_t$ , indexed by  $P$ , along with linear maps  $u_t^s : U_s \rightarrow U_t$  for  $s \leq t$ . This module can be viewed as a functor from  $P$  to the category of vector spaces or  $R$ -modules. Persistence modules over a field  $F$  indexed by  $\mathbb{N}$  can be expressed as:

$$U \simeq \bigoplus_i x^{t_i} \cdot F[x] \oplus \left( \bigoplus_j x^{r_j} \cdot (F[x]/(x^{s_j} \cdot F[x])) \right) \quad (27)$$

Here, multiplication by  $x$  represents a forward step in the persistence module. The free parts correspond to homology generators that appear at a certain filtration level and persist indefinitely, whereas torsion parts correspond to those that appear at a filtration level and last for a finite number of steps (Barannikov, 1994; Zomorodian and Carlsson, 2004).

This framework allows the unique representation of the persistent homology of a filtered simplicial complex using either a persistence barcode or a persistence diagram. In the barcode, each persistent

generator is represented by a line segment starting and ending at specific filtration levels, while in the diagram, each generator is represented as a point with coordinates indicating its birth and death times. Barannikov’s canonical form offers an equivalent representation.

## C.2 Stability

The stability of persistent homology is a key attribute, particularly in its application to data analysis, as it ensures robustness against small perturbations or noise in the data (Cohen-Steiner et al., 2005). This stability is quantitatively defined in terms of the **bottleneck distance**, a metric for comparing persistence diagrams.

The bottleneck distance between two persistence diagrams  $X$  and  $Y$  is defined as:

$$W_\infty(X, Y) := \inf_{\varphi: X \rightarrow Y} \sup_{x \in X} \|x - \varphi(x)\|_\infty \quad (28)$$

where the infimum is taken over all bijections  $\varphi$  from  $X$  to  $Y$ . This metric essentially measures the greatest distance between matched points (or generators) in two persistence diagrams, considering the optimal matching.

A fundamental result in the theory of persistent homology is that small changes in the input data (such as a filtration of a space) result in small changes in the corresponding persistence diagram, as measured by the bottleneck distance. This is formalized by considering a space  $X$ , homeomorphic to a simplicial complex, with a filtration determined by the sublevel sets of a continuous tame function  $f : X \rightarrow \mathbb{R}$ . The map  $D$  that takes the function  $f$  to the persistence diagram of its  $k$ th homology is 1-Lipschitz with respect to the supremum norm on functions and the bottleneck distance on persistence diagrams. Formally, this is expressed as (Cohen-Steiner et al., 2005):

$$W_\infty(D(f), D(g)) \leq \|f - g\|_\infty \quad (29)$$

This Lipschitz condition implies that a small change in the function  $f$ , as measured by the supremum norm, will not cause a disproportionately large change in the persistence diagram. Consequently, persistent homology is particularly useful in applications where data may be subject to noise or small variations, as the essential topological features (captured by the persistence diagrams) are not overly sensitive to such perturbations.

## C.3 Computation

There are various software packages for computing persistence intervals of a finite filtration (Otter et al., 2017). The principal algorithm is based on the bringing of the filtered complex to its canonical form by upper-triangular matrices (Barannikov, 1994).

## D Complete Experimental Results of Neurological Topology Clustering

In Section 3, it is mentioned that the lengths of DKNs vary. For ease of presentation, we have only shown cases with a larger amount of data. Here, we provide the complete table, as shown in Table 4, which is generated under the same settings as used for Table 1.

For our method, we have also created line charts for DKNs of each length, serving as a supplement to Figure 4. We use the same background color as the figures in the main text to represent images that are completely consistent with the main text, while a white background is used for the line charts corresponding to DKNs of other lengths. Figure 12 show the results corresponding to DKN obtained using the NTC method. Additionally, we also present the results corresponding to other methods. Figure 13 display the results corresponding to DKN obtained using the DBSCAN method. Figure 14 show the results corresponding to DKN obtained using the Hierarchical method. Figures 15 and 16 illustrate the results corresponding to DKN obtained using the K-Means method. Since the DKN length for the method corresponding to AMIG (Chen et al., 2024a) is limited to 2, there is no need to display the inflection points.

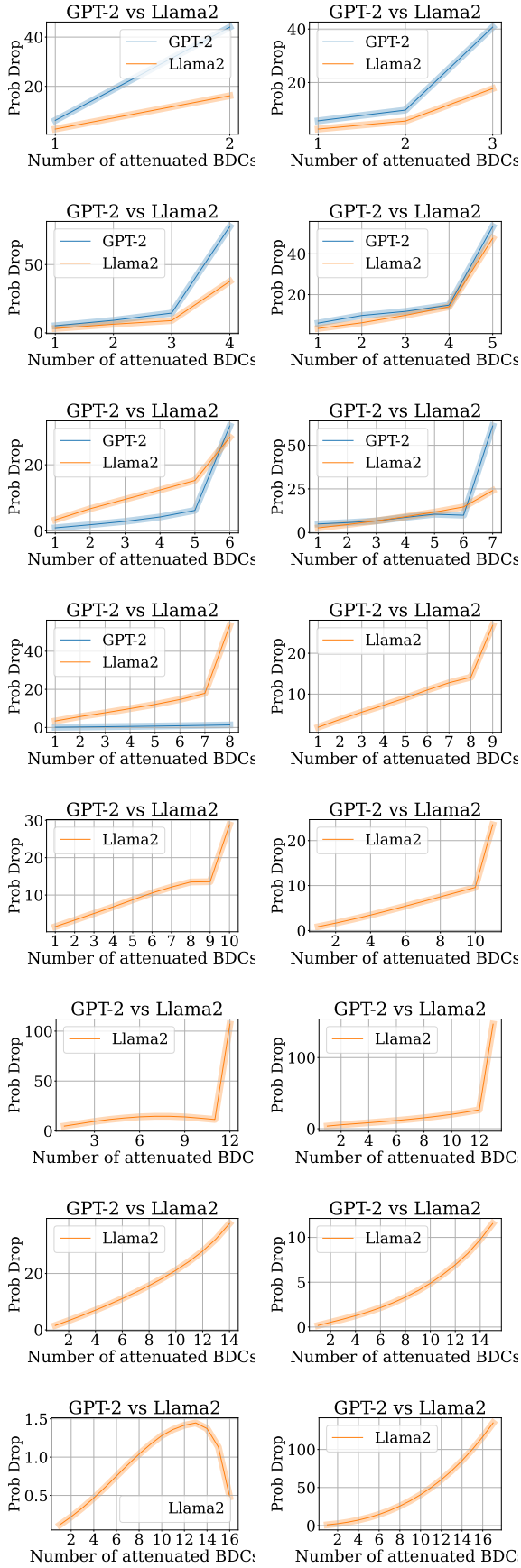


Figure 12: The graph of  $\Delta Prob$  relative to the number of suppressed BDCs obtained using the NTC method.

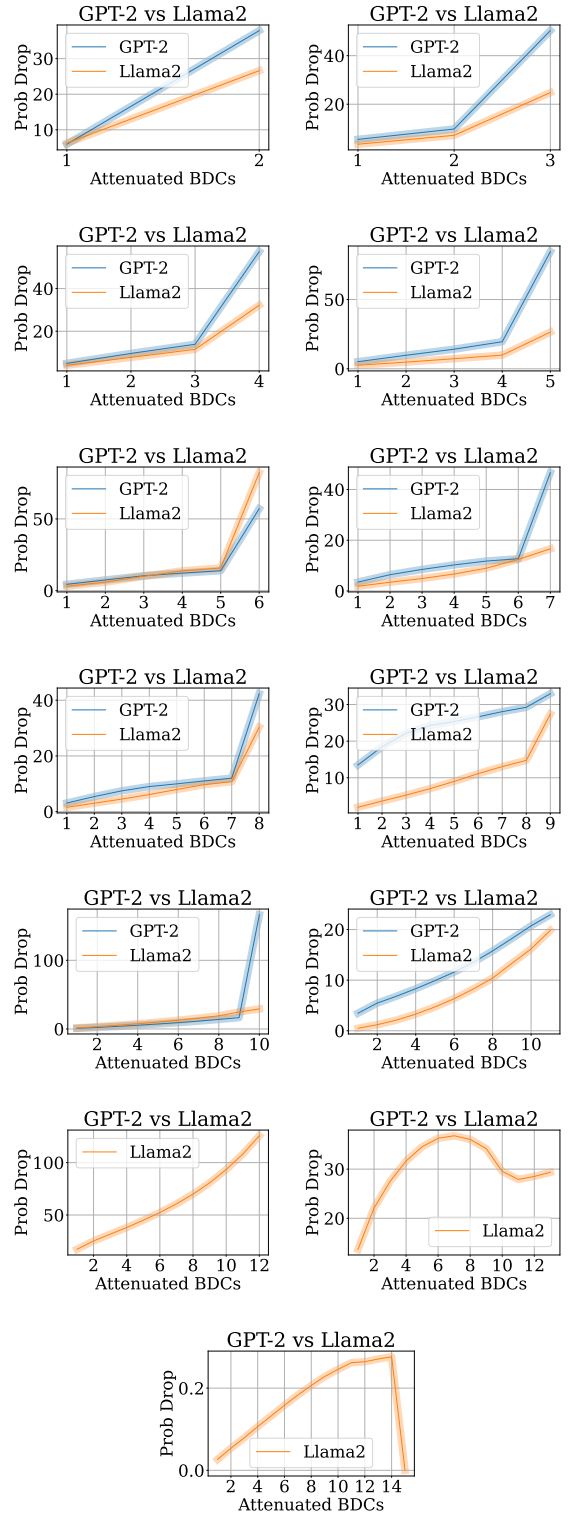


Figure 13: The graph of  $\Delta Prob$  relative to the number of suppressed BDCs obtained using the DBSCAN method.



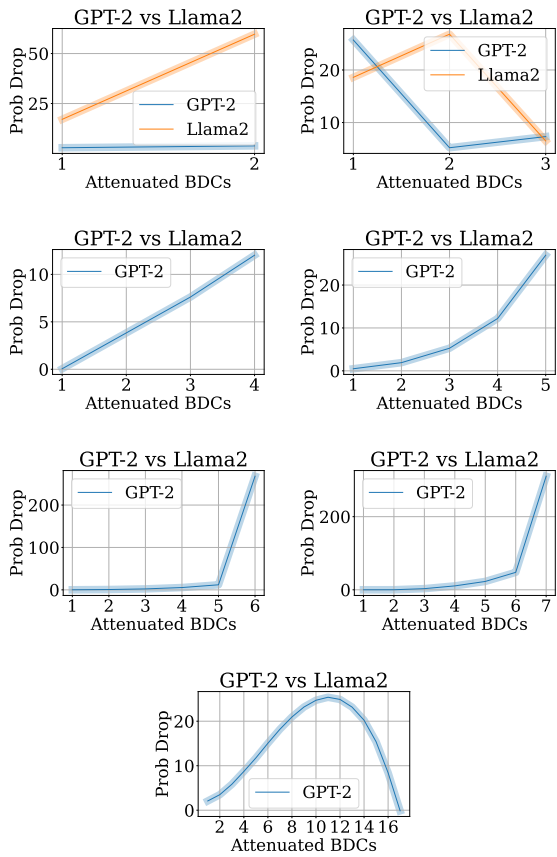


Figure 14: The graph of  $\Delta Prob$  relative to the number of suppressed BDCs obtained using the Hierarchical method.

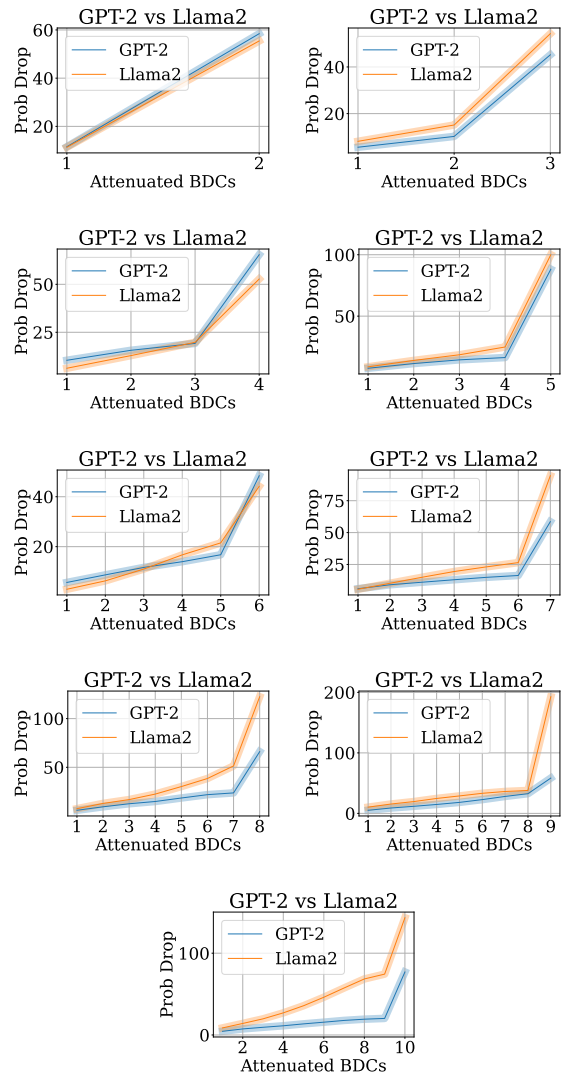


Figure 15: The graph of  $\Delta Prob$  relative to the number of suppressed BDCs obtained using the K-Means method (Part 1).

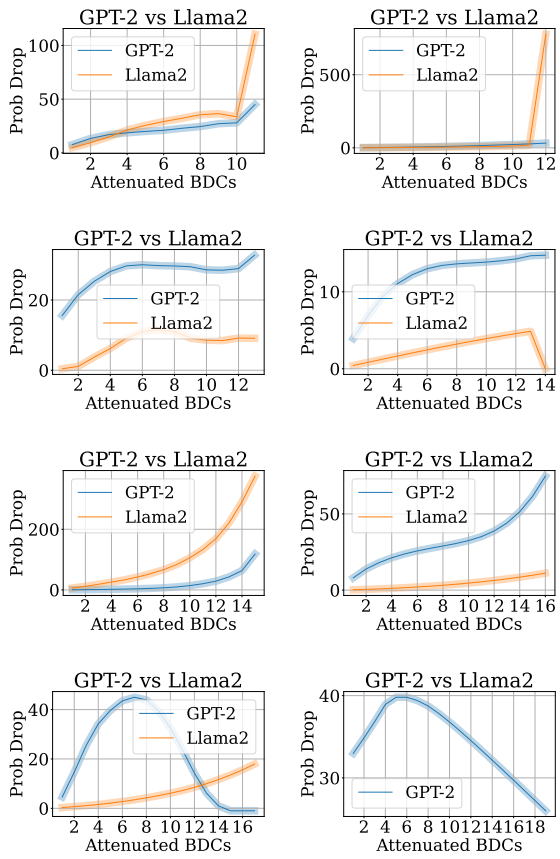


Figure 16: The graph of  $\Delta Prob$  relative to the number of suppressed BDCs obtained using the K-Means method (Part 2).

Method	GPT-2						
	1	2	3	4	5	6	7
DBSCAN	—	12 → 26	17.7 → 38.8	23.4 → 33.8	32.9 → 53.4	24.8 → 49.8	26.4 → 46.6
Hierarchical	—	2.9 → 3.8	15.5 → 7.3	3.7 → 12	4.3 → 27	25 → 92	—
K-means	—	27.4 → 38.3	17.8 → 29.5	27.8 → 44.6	32.4 → 58.7	28.9 → 30.5	21.8 → 30.1
AMIG	—	-0.79 → 0.89	—	—	—	—	—
<b>NTC (Ours)</b>	—	<b>7.9 → 53</b>	<b>8.6 → 44</b>	<b>12 → 92</b>	<b>11 → 53</b>	<b>3.1 → 32</b>	<b>7.8 → 61</b>
Method	GPT-2						
	8	9	10	11	12	13	14
DBSCAN	9.1 → 12.5	30.4 → 32.9	39.8 → 19.3	11 → 22.9	—	—	—
Hierarchical	—	—	—	—	—	—	—
K-means	23.3 → 33.6	25.1 → 58.4	18.5 → 47.3	31 → 45	10.3 → 32.1	40.6 → 32.7	13 → 14.8
AMIG	—	—	—	—	—	—	—
<b>NTC (Ours)</b>	<b>0.6 → 1.3</b>	—	—	—	—	—	—
Method	GPT-2						
	15	16	17	18	19	Overall	
DBSCAN	—	—	—	—	—	23.90 → 34.72	
Hierarchical	—	—	20.8 → -0.13	—	—	20.77 → 25.05	
K-means	48.8 → 118.1	29.2 → 74.7	62.7 → -1	—	37 → 26	34.42 → 38.86	
AMIG	—	—	—	—	—	-0.79 → 0.89	
<b>NTC (Ours)</b>	—	—	—	—	—	<b>9.32 → 55.60</b>	
Method	Llama2						
	1	2	3	4	5	6	
DBSCAN	—	7.1 → 15.4	8.7 → 15.2	11.2 → 21.0	8.9 → 14.6	17.6 → 34.2	
Hierarchical	—	27.6 → 44.3	22.3 → 6.5	—	—	—	
K-means	—	2.8 → 16.1	19.2 → 39.1	22.2 → 45.4	25.8 → 41.4	14.2 → 34.2	
AMIG	—	-0.99 → 1.99	—	—	—	—	
<b>NTC (Ours)</b>	—	<b>2.8 → 15.6</b>	<b>4.3 → 19.1</b>	<b>6.4 → 37.4</b>	<b>8.4 → 49.9</b>	<b>9.8 → 29.6</b>	
Method	Llama2						
	7	8	9	10	11	12	
DBSCAN	7.7 → 18.0	7.7 → 16.3	12.4 → 28.6	12.2 → 32.1	7.2 → 25.3	53.6 → 125.6	
Hierarchical	—	—	—	—	—	—	
K-means	28.6 → 50.8	37.9 → 97.1	29.4 → 74.9	47.9 → 87.8	50.7 → 110.5	13.5 → 23.2	
AMIG	—	—	—	—	—	—	
<b>NTC (Ours)</b>	<b>9.1 → 27.6</b>	<b>7.8 → 50.1</b>	<b>7.6 → 21.5</b>	<b>9.1 → 31.5</b>	<b>4.6 → 25.7</b>	<b>15.2 → 125.5</b>	
Method	Llama2						
	13	14	15	16	17	Overall	
DBSCAN	92.5 → 29.4	—	—	0.2 → 0	—	22.26 → 18.24	
Hierarchical	9.7 → 9.1	36.3 → 20.2	—	—	—	27.50 → 44.04	
K-means	19.7 → 19.2	36.3 → 50.0	68.8 → 375.0	3.3 → 11.1	4.9 → 17.9	20.08 → 39.24	
AMIG	—	—	—	—	—	-0.99 → 1.99	
<b>NTC (Ours)</b>	<b>15.8 → 183.7</b>	<b>13.6 → 37.7</b>	<b>3.1 → 17.2</b>	<b>1.0 → 0.5</b>	<b>31.2 → 134.9</b>	<b>14.11 → 28.78</b>	

Table 4: The full results of Comparison of the degeneracy properties of DKN using different methods. “—” denotes methods that do not yield  $\mathcal{D}$  of a specific length.