# Knowledge Homophily in Large Language Models

**Utkarsh Sahu**[1], **Zhisheng Qi**[1], **Mahantesh Halappanavar**[2], **Nedim Lipka**[3], **Ryan A. Rossi**[3],
**Franck Dernoncourt**[3], **Yu Zhang**[4], **Yao Ma**[5], **Yu Wang**[1]

[1]University of Oregon  [2]Pacific Northwest National Laboratory  [3]Adobe Research  [4]Texas A&M University  [5]Rensselaer Polytechnic Institute

{utkarsh,charq,yuwang}@uoregon.edu, hala@pnnl.gov, {lipka,ryrossi,dernonco}@adobe.com, yuzhang@tamu.edu, may13@rpi.edu

## ABSTRACT

Large Language Models (LLMs) have been increasingly studied as neural knowledge bases for supporting knowledge-intensive applications such as question answering and fact checking. However, the structural organization of their knowledge remains unexplored. Inspired by cognitive neuroscience findings, such as semantic clustering and priming, where knowing one fact increases the likelihood of recalling related facts, we investigate an analogous knowledge homophily pattern in LLMs. To this end, we map LLM knowledge into a graph representation through knowledge checking at both the triplet and entity levels. After that, we analyze the knowledge-ability relationship between an entity and its neighbors, discovering that LLMs tend to possess a similar level of knowledge about entities positioned closer in the graph. Motivated by this homophily principle, we propose a Graph Neural Network (GNN) regression model to estimate entity-level knowledgeability scores for triplets by leveraging their neighborhood scores. The predicted knowledge-ability enables us to prioritize checking less well-known triplets, thereby maximizing knowledge coverage under the same labeling budget. This not only improves the efficiency of active labeling for fine-tuning to inject knowledge into LLMs but also enhances multi-hop path retrieval in reasoning-intensive question answering. Our code is available at https://github.com/utkarshxsahu/kgc.
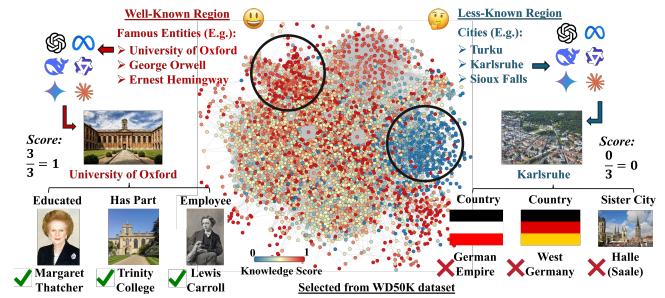
## KEYWORDS

Structured Knowledge, Large Language Model, Homophily

## 1 INTRODUCTION

Large Language Models (LLMs) have emerged as powerful neural knowledge bases by encoding vast amounts of world knowledge within their neural parameters [10, 21]. This neural-embedded knowledge enables LLMs to produce contextually relevant and factually rich responses, supporting real-world applications such as fact checking [11] and question answering [32, 34]. To better explore this neural knowledge base, researchers have devised knowledge checking methods to investigate the knowledge patterns of LLMs [1, 38] and leveraged the derived insights to guide knowledge-intensive tasks, including adaptive retrieval [35, 36], knowledge editing [25], and hallucination detection [26].

Despite various knowledge patterns identified previously [10, 21, 37], little attention has been given to whether LLMs' knowledge exhibits structural organization. In fact, in cognitive neuroscience [12], several works have highlighted the semantic clustered patterns of the neural knowledge in human brain networks: (i) semantic clustering in memory recall, where people tend to retrieve related words together (e.g., recalling "dog, cat, horse" in sequence) [2, 16], and (ii)



**Figure 1: LLM is prompted on individual triple facts , which are aggregated to obtain entity-level knowledgeability scores. When visualized, these scores reveal the knowledge homophily pattern, where topologically close entities form distinct high-knowledge (red) and low-knowledge (blue) communities. Note that graph layout is visualized by the ForceAtlas2 algorithm [8] to preserve topological proximity.**

homophily brain networks, where regions with similar functions or inputs are more likely to connect [28]. Analogously, we hypothesize that LLMs also exhibit a similar knowledge homophily pattern, i.e., they tend to possess similar levels of knowledge about conceptually related entities, as illustrated in Figure 1 by checking GPT-3.5's knowledge about triplets from WD50K dataset. Discovering this phenomenon provides valuable insights into the structural knowledge organization in LLMs and the design of knowledge-intensive task solutions. For example, by estimating the knowledgeability of concepts based on related concepts, one can identify factually weaker areas, thereby enabling more efficient knowledge labeling for knowledge injection and retrieval as we conducted in Section 4.

Given the existence of knowledge homophily in other disciplines and its potential implications, this paper uncovers this innovative pattern and designs graph-based machine learning models to exploit neighborhood information to predict knowledgeability scores and discover less-known regions. These insights then guide efficient labeling for LLM fine-tuning and improve knowledge retrieval for multi-hop reasoning in question answering. Our contributions are:

- **Knowledge Homophily Discovery:** We demonstrate the existence of knowledge homophily in LLMs by measuring knowledge at triplet/entity levels, showing that topologically close entities tend to exhibit similar knowledgeability scores.
- **Knowledge Homophily Application:** We leverage the discovered knowledge homophily to develop a GNN-based estimator that infers the entity knowledgeability, and showcase two applications enhancing knowledge injection efficiency and guiding multi-hop retrieval for question-answering.

## 2 RELATED WORK

**Knowledge Checking for LLMs as Knowledge Bases (KBs).**
As model parameters keep scaling, LLMs have evolved from task-specific executors to general-purpose agents and neural knowledge bases capable of storing and providing factual information for knowledge-intensive applications [20, 23]. However, unlike traditional KBs with structured schemas [30], the knowledge stored in LLM is encoded implicitly and non-interpretably within its parameters. This inherent lack of transparency creates a critical need for methods to verify what these models "know" and to ensure their reliable utilization, motivating numerous work on knowledge checking. These methods can be broadly categorized by their goals: verifying factual accuracy [7, 11], assessing self-awareness [10, 29], and gauging knowledge coverage and consistency of LLMs against their own internal knowledge and external knowledge in other formats [13, 15]. While valuable, they primarily focus on the content perspective of LLMs' internal knowledge instead of the structural perspective reviewed next.

**Pattern Understanding of LLMs as Knowledge Bases.** With LLMs increasingly used as neural knowledge bases, existing studies have sought to uncover the underlying structured knowledge patterns encoded in their parameters from two perspectives. The first group investigates where knowledge is stored within the neural parameters of LLMs. For example, studies have shown that feed-forward layers can act as key-value memories for specific facts [6] and that factual associations are often localized and editable within these mid-layer modules, sometimes referred to as "knowledge neurons" [3, 17]. The second focuses on how knowledge is structurally organized. [19] evaluated LLMs on structural properties like symmetry, hierarchy, and path-following, showing that models fail on complex relational tests. These structural analyses of knowledge patterns suggest the existence of implicit structures within LLMs, yet the homophily property remains underexplored.

## 3 KNOWLEDGE HOMOPHILY DISCOVERY

In this section, we compute triplet-level knowledgeability scores and aggregate these measures to obtain entity-level knowledgeability scores. After that, we assess homophily by quantifying the knowledgeability differences between neighboring entities in Section 3.2.1 and qualitatively visualizing their scores in Section 3.2.2, thereby revealing the knowledge homophily phenomenon.

### 3.1 Knowledgeability Computation

To examine whether LLMs exhibit consistent knowledge about neighboring entities, we first evaluate knowledgeability at the triplet level and then aggregate it to obtain an entity-level score. Given triplets $\mathcal{T} = \{(s_i, d_i, r_i)\}_{i=1}^{|\mathcal{T}|}$ from the knowledge graph, where a source entity $s_i$ is connected to a destination entity $d_i$ via relation $r_i$, we define the knowledgeability of the LLM on triplet $(s_i, d_i, r_i)$ as $\mathcal{K}(s_i, d_i, r_i)$, reflecting how well the LLM knows about this relational fact. For each entity $s_i$, we denote its neighbor entity set as $\mathcal{N}(s_i)$, representing the entities adjacent to $s_i$ as either head or tail, and their corresponding neighbor triplet set as $\mathcal{T}(s_i)$. The entity-level knowledgeability of $s_i$, denoted as $\mathcal{K}(s_i)$, is derived by aggregating knowledgeability scores over its neighboring triplets, capturing how well the LLM knows about entity $s_i$. Next, we introduce details of calculating triplet and entity knowledgeability.

---

> **Prompt 1: LLM-based Triplet Evaluation**
>
> **System Message:** Evaluate the statement based on your knowledge and respond with `True` or `False`.
>
> **Given:** Triplet $\mathcal{T} = (sub, rel, obj)$, **Date $D$ (Temporal Version)**
>
> **Template:** Relation → pattern (e.g., son_of → {SUB} is the son of {OBJ}.)
>
> **Procedure:**
> (1) Retrieve relation-based template for *rel* in triplet $\mathcal{T}$.
> (2) Fill {SUB}→*sub*, {OBJ}→*obj* from $\mathcal{T}$ to get statement $\mathcal{S}$.
> (3) **Append $\mathcal{S}$ on Date $D$ (Temporal Version)**
> (4) Prompt **System Msg** + **User Msg:** $\mathcal{S}$ to the LLM.
> (5) Record the model output in the format of `True/False`

*3.1.1 Calculating Triplet Knowledgeability.* Following prior work showing that LLMs are generally well-calibrated in knowing what they know [1, 10, 21], we convert each triplet $(s_i, d_i, r_i)$ into a natural language statement and prompt the LLM to judge whether it recognizes the fact. The model response is recorded as a binary value, with `True`/`False` mapping to 1/0, representing its knowledgeability about the triplet $\mathcal{K}(s_i, d_i, r_i)$. For temporal triplets $(s_i, d_i, r_i, t)$ (e.g., "Donald Trump made a visit to China on 2017-11-08."), we extend the prompt to include the timestamp, allowing us to assess the temporal dimension of LLM knowledgeability. Prompt 1 illustrates the template, with temporal variants highlighted in red.

*3.1.2 Calculating Entity Knowledgeability.* Given the above calculated triplet knowledgeability, we obtain knowledgeability of the entity $v_i$ by aggregating from all triplets involving $v_i$ [9, 22]:

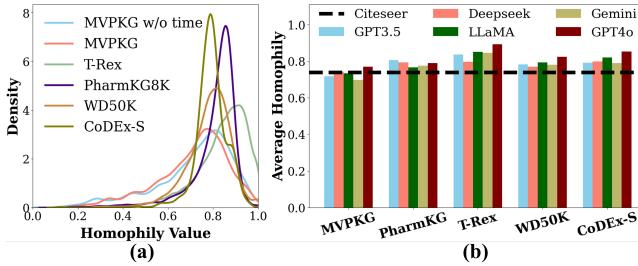$$\mathcal{K}(v_i) = |\mathcal{T}(v_i)|^{-1} \sum_{(s,d,r) \in \mathcal{T}(v_i)} \mathcal{K}(s, d, r). \quad (1)$$

Note that the above neighborhood aggregation naturally extends to temporal triplets $(s, d, r, t) \in \mathcal{T}(v_i)$, allowing temporal information to be incorporated into the entity knowledgeability calculation.

### 3.2 Homophily Computation and Analysis

Furthermore, we evaluate whether topologically close entities share similar knowledgeability, i.e., the homophily of entity knowledgeability $\mathcal{H}(v_i)$. Inspired by existing homophily computation [14, 31], we compute knowledgeability homophily as one minus the average absolute difference in knowledgeability between central node $v_i$ and its neighbors $\mathcal{N}(v_j)$ in the knowledge graph:

$$\mathcal{H}(v_i) = 1 - \frac{1}{|\mathcal{N}(v_i)|} \sum_{v_j \in \mathcal{N}(v_i)} |\mathcal{K}(v_i) - \mathcal{K}(v_j)| \quad (2)$$

where a smaller difference between neighboring entities, $|\mathcal{K}(v_i) - \mathcal{K}(v_j)|$, leads to a higher homophily value $\mathcal{H}(v_i)$. We empirically quantify triplet and entity knowledgeability scores and analyze homophily patterns both quantitatively and qualitatively. We evaluate five representative LLMs: GPT-3.5, 4o, Gemini-2.5 Flash, LLaMA3.3-70B, and DeepSeek-V3. These models are assessed across five knowledge graphs: MVPKG [18], T-Rex [4], PharmKG [39], WD50K [5], and CoDEx-S [24]. Among them, T-Rex, WD50K, and CoDEx-S represent general Wikipedia knowledge, whereas PharmKG8K and MVPKG focus on biomedical and political science. Our graph visualizations in Figure 1/3 employ the ForceAtlas2 [8] to position topologically close nodes visually close, enabling an intuitive assessment of whether they share similar knowledgeability scores.
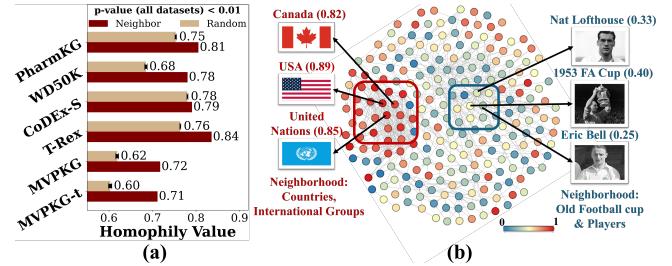
Figure 2: (a): Distribution of node knowledgeability homophily for each dataset; (b): Average knowledge homophily across datasets and LLMs with black line showing Citeseer benchmark (0.74) for high homophily in graph analysis.

*3.2.1 Quantitative Analysis of Knowledge Homophily.* Figure 2(a)/(b) presents the node homophily distribution and the average graph homophily across several knowledge graphs. In Figure 2(a), these distributions are all right-skewed and peak around 0.8, suggesting the majority portion of nodes and their neighbors tend to share similar knowledgeability scores. This high homophily property enhances graph machine learning in node-level prediction, such as node classification and regression [40], and therefore inspires regression to predict entities' knowledge scores in Section 4.1. Furthermore, incorporating temporal information into MVPKG results in a slight left shift. This slightly decreasing neighbor homophily score might indicate that the temporal dimension introduces greater complexity and finer distinctions in knowledgeability between the nodes and their neighbors. In addition, we compute the average graph homophily by averaging across all nodes in Figure 2(b). Across datasets and LLMs, the scores mostly exceed the homophily of the Citeseer [27, 33] benchmark (horizontal line), a classic homophily network for node classification. This consistency across diverse settings indicates that the observed homophily generally aligns with the conventional notion of a "high-homophily" network.

Furthermore, we compare the calculated knowledge homophily against a degree-matched random baseline. For each node $v$ instead of using its actual neighborhood $\mathcal{N}(v)$, we create a random "peer group" $\widehat{\mathcal{N}}(v)$ by randomly sampling the same number of entities $(|\widehat{\mathcal{N}}(v)| = |\mathcal{N}(v)|)$ from the entire graph. Then, a homophily score is calculated as the difference between the central node $v$ and this randomly sampled group. The true neighborhood homophily was found to be statistically significant than this random baseline (100 trials per dataset, two-tailed z-test, $p < 0.01$), with extremely narrow 99% confidence intervals. *The results in Figure 3(a) show that across datasets, the homophily calculated from true neighbors exceeds the random baseline, confirming that it is not a statistical artifact or random fluctuation, but rather a robust and intrinsic structural property of LLMs' internal knowledge organization.*

*3.2.2 Qualitative Analysis of Knowledge Homophily.* Figure 3(b) shows a sampled T-Rex subgraph with nodes colored by entity-level knowledgeability $\mathcal{K}(v)$. A geopolitical neighborhood (countries) forms a compact and uniformly high-knowledgeability region with small intra-neighborhood deviations, confirming the high knowledgeability homophily. A historical football neighborhood (event/players) is similarly coherent but at lower knowledgeability. Despite varying means, both exhibit knowledge homophily.



Figure 3: (a) Neighboring nodes possess similar knowledgeability scores than randomly sampled nodes, indicating that homophily is a intrinsic structural property of LLMs; (b) Entities with their distinct knowledgeability levels $\mathcal{K}(v)$ indicated by node color (**Red** = High, **Blue** = Low).

## 4 KNOWLEDGE HOMOPHILY APPLICATION

After identifying the knowledge homophily pattern in Section 3, where topologically proximate entities exhibit similar knowledgeability scores, we apply this finding to two knowledge-intensive tasks: (1) homophily-aware knowledge checking for more efficient fine-tuning, and (2) homophily-aware knowledge retrieval for enhanced question answering. The core idea is to train a GNN-based model to estimate entity-level knowledgeability based on their neighbor knowledgeability and thereby identify triplets in less-known regions. These triplets can then be prioritized either for knowledge checking, to maximize fine-tuning benefits when injecting knowledge into LLMs, or for retrieval, to complement missing knowledge of the downstream answer generator. Both tasks rely on knowledgeability estimation to pinpoint knowledge gaps.

## 4.1 Homophily-aware Knowledge Estimation

Given homophily is sufficient for high-utility GNN predictions [14], we design a GNN-based regression model to perform message-passing, aggregate neighboring entity embeddings, and predict previously unknown entity scores. Specifically, given a set of entities $\mathcal{V}^{\text{Train}}$ with known knowledgeability (by prompting LLMs), our goal is to train a GNN to estimate the knowledgeability of unseen entities. At each layer, the model performs Message Passing (MP) and Feature Transformation (TR), followed by regression:

$$\widehat{\mathcal{K}}_i^l = \text{MP}^l\left(\widetilde{\mathcal{K}}_j^{l-1} \mid v_j \in \mathcal{N}(v_i) \cup v_i\right), \quad \widetilde{\mathcal{K}}_i^l = \text{TR}^l(\widehat{\mathcal{K}}_i^l), \qquad (3)$$

$$\mathcal{L} = \frac{1}{|\mathcal{V}^{\text{Train}}|} \sum_{v_i \in \mathcal{V}^{\text{Train}}} \left|\widetilde{\mathcal{K}}_i^l - \mathcal{K}_i\right|^2, \qquad (4)$$

The initial node feature matrix is defined as $\widetilde{\mathcal{K}}^0 = [\mathcal{X}_1, \ldots, \mathcal{X}_{v_{|\mathcal{V}|}}]^\top$, where each node feature $\mathcal{X}_{v_i}$ is a dense text embedding obtained from pretrained language models. After training on $\mathcal{V}^{\text{Train}}$, the model is further used to infer the knowledgeability scores for all entities in the knowledge graph, eliminating the need for resource/time-intensive knowledge probing via exhaustive LLM prompting. In the following, we will utilize these estimated knowledgeability scores of entities to either guide triplet selection for subsequent LLM fine-tuning in Section 4.2 or guide retrieval for answering reasoning-intensive multi-hop questions in Section 4.3. Due to space constraints, we only briefly describe the experimental setup and provide full details in Figure 4 in Appendix.

**Table 1: Performance comparison of fine-tuning with triplets selected by knowledgeability estimated by Random, MLP, and GNN. Best result in bold and second-best underlined. L=Llama3-8B, M=Mistral-7B. Selection Quality: percentage of triplets selected for fine-tuning that are unknown to base LLMs. Generalization Gain: percentage of additional 2% evaluation triplets identified by the fine-tuned LLMs. Detailed setting is visualized in Figure 4 in Appendix.**

| Task | Method | T-Rex | | PharmKG | | WD50K | | MVPKG | | CoDExS | | Avg. |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | L | M | L | M | L | M | L | M | L | M | |
| **Selection Quality** | Rand | 36.5 | 44.8 | 81.9 | 72.8 | 41.2 | 46.8 | <u>68.5</u> | 66.3 | 33.8 | 51.7 | 54.4 |
| | MLP | **38.4** | <u>48.7</u> | <u>84.8</u> | <u>77.0</u> | <u>44.3</u> | <u>47.8</u> | 67.9 | <u>68.6</u> | <u>39.1</u> | <u>57.8</u> | <u>57.4</u> |
| | GNN | <u>37.3</u> | **54.5** | **87.6** | **79.2** | **49.6** | **50.6** | **72.2** | **71.5** | **45.5** | **63.9** | **61.2** |
| **Generalization Gain** | Base | 63.3 | 64.0 | 17.8 | 55.3 | 54.8 | 42.9 | 26.1 | 52.3 | 64.9 | 58.5 | 49.9 |
| | Rand | 86.4 | 81.9 | 34.9 | 41.3 | <u>57.8</u> | **56.3** | 30.7 | 65.1 | **78.8** | 72.1 | 60.5 |
| | MLP | <u>87.9</u> | <u>90.2</u> | <u>35.8</u> | <u>57.2</u> | 56.1 | 53.2 | <u>42.8</u> | <u>74.5</u> | 73.7 | <u>85.2</u> | <u>65.6</u> |
| | GNN | **89.1** | **91.9** | **37.0** | **60.7** | **58.8** | <u>55.1</u> | **44.5** | **76.7** | <u>75.6</u> | **88.0** | **67.7** |

## 4.2 Homophily-guided Knowledge Injection

We leverage the homophily to estimate triplet knowledgeability and prioritize selecting less-known triplets for fine-tuning LLMs within a fixed budget, thereby enabling more effective knowledge injection into LLMs. For each dataset, we allocate 4000 triplets as the knowledge-checking budget for selection and fine-tuning, with an additional 2% of all triplets reserved as the test set. Within the 4000 budget, 20% of triplets are sampled as anchor points, for which we directly query the base LLM to obtain ground-truth binary knowledgeability scores (Section 3.1). These anchors provide labeled data to estimate the knowledgeability of their associated entities, which is used to train a GNN model (Eq. (4)) and predict knowledgeability scores for all remaining entities. Based on these predictions, we prioritize triplets with lower-scored entities from the remaining 80% unqueried pool to complete the 4000-triplet set for fine-tuning. We benchmark this triplet selection against two baselines: Random, which uniformly samples triplets, and MLP, which estimates knowledgeability without homophily, eliciting the knowledge homophily contribution to knowledge estimation. We experiment with LLaMA3-8B(L) and Mistral-7B(M).

Table 1 evaluates homophily-guided knowledge injection from two perspectives: selection quality and generalization gain. For selection quality, we assess whether the chosen triplets better capture the knowledge deficiencies of LLMs. Among the 4000 triplets selected for fine-tuning, we compute the percentage that the base LLM does not recognize, following the procedure in Section 3.1. A higher score indicates that more selected triplets are unknown to the LLM and thus more valuable for fine-tuning. Our GNN regressor achieves the highest proportion of unknown triplets, outperforming MLP and Random selection. This demonstrates the advantage of incorporating homophily into GNN design in enabling more effective estimation of ground-truth knowledgeability for knowledge injection. For generalization gain, we test whether fine-tuning on selected triplets improves the knowledgeability over the reserved 2% held-out set. The best performing GNN regressor confirms that its higher selection precision translates into stronger knowledge generalization. This superior generalization gain holds across different evaluation budgets from 1% to 20% in Figure 5 in Appendix.

**Table 2: Multi-hop Question Answering Accuracy by GPT4-as-a-Judge; M=MLP, G=GNN, BS=Beam Search, H = Hop**

| Dataset<br>Q-Hop | T-Rex | | PharmKG | | WD50K | | MVPKG | | CoDExS | |
|---|---|---|---|---|---|---|---|---|---|---|
| | 2-H | 3-H | 2-H | 3-H | 2-H | 3-H | 2-H | 3-H | 2-H | 3-H |
| **Base** | 30.9 | 22.6 | 21.4 | 16.0 | 25.1 | 17.2 | 24.4 | 19.1 | 28.6 | 20.4 |
| **M-BS** | 33.8 | 23.1 | 21.7 | 16.2 | 25.8 | 17.3 | 24.9 | 19.4 | 29.9 | 20.5 |
| **G-BS** | **34.2** | **23.7** | **22.2** | **16.6** | **26.0** | **17.5** | **25.4** | **19.5** | **31.1** | **20.9** |

## 4.3 Homophily-guided Knowledge Retrieval

We test whether the estimated knowledgeability can guide entity retrieval to provide better context for question-answering. For each KG, we generate 1000 questions (500 2-hop/500 3-hop). Entity knowledgeability $\mathcal{K}(v)$ is predicted with a GNN regressor trained on 40% of entities labeled by GPT-3.5, excluding entities for generating 1000 evaluation questions. We embed both entities/relations and questions using `all-MiniLM-L6-v2`. Starting with entity linking in the question, we run beam search up to the hop limit and score each neighbor by its knowledgeability $\mathcal{K}(v)$ and semantic similarity $\mathcal{S}(r||d, q)$ to the question $q$ where $r||d$ represents its relation $r$ concatenated with the tail entity $d$. Baselines are as follows:

- **Baseline (Semantic Beam Search):** It retrieves paths using beam search guided solely by the semantic similarity $\mathcal{S}(r||d, q)$ between the path (relation + tail entity) and the input question.
- **Knowledge-aware Beam Search (BS):** This method adjusts beam search to favor less-known entities. For each expansion, the semantic score $\mathcal{S}$ is penalized by the next entity knowledgeability, $\mathcal{K}(u)$. The final score is $\mathcal{S} \times (1 - \alpha \cdot \mathcal{K}(u))$ with $\alpha$ being weighting factor. Entities with lower knowledgeability receive higher retrieval priority, achieving knowledge-aware search. Beam Search (BS) with GNN/MLP as knowledge estimator are G-BS/M-BS.

The reader model GPT-3.5 generates answers strictly based on the retrieved triples, and their correctness is evaluated by GPT-4. In Table 2, M/G-BS outperforms Baseline across all datasets and hops. Crucially, GNN-based beam search (G-BS) surpasses its homophily-agnostic MLP-based counterpart (M-BS), confirming the advantage of structural knowledge homophily in estimating knowledgeability. On 2-hop, G-BS has an improvement of 4.57% over the baseline, with larger gains on general KGs (e.g., T-Rex) than domain-specific ones (e.g., PharmKG). For 3-hop questions, all methods drop due to longer chains and semantic drift, but G-BS still has an improvement of 2.62% over the baseline. This superior performance of G-BS holds across different budgets for training the knowledgeability estimator in Figure 6 in Appendix.
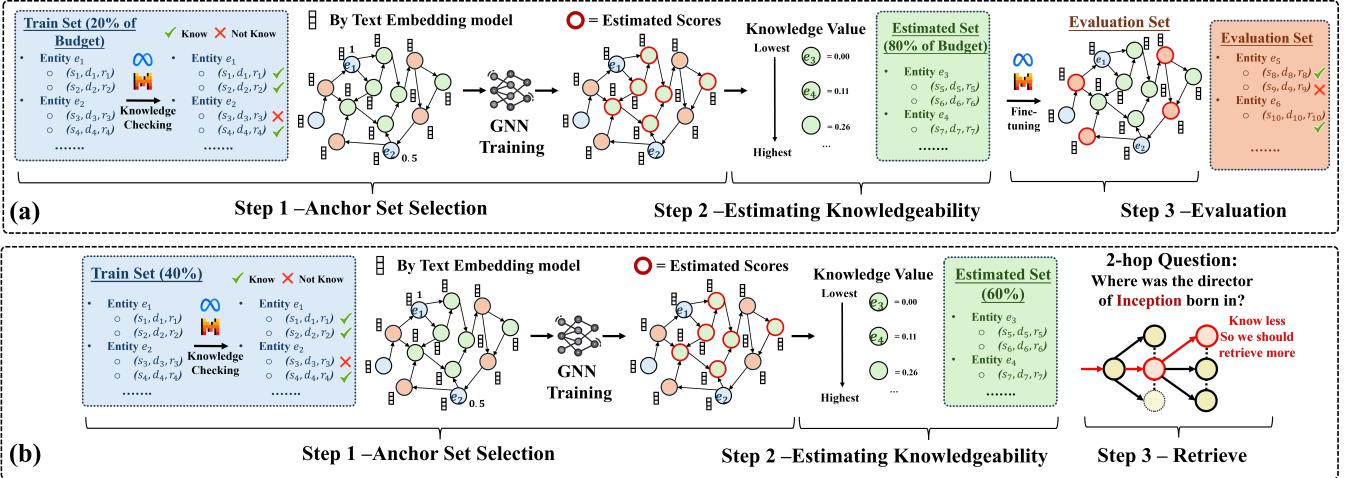
## 5 CONCLUSION

Motivated by the structural knowledge organization of the human brain, this work explores the homophily structure of neural knowledge stored in LLMs, showing that LLMs' knowledge about entities is strongly correlated with their neighbors in a knowledge graph. Building on this insight, we design a GNN regressor to estimate entity-level knowledgeability scores by leveraging their local neighborhood scores. We demonstrate the utility of these scores in two applications: selecting less-known triplets for efficient knowledge injection through fine-tuning, and improving retrieved context for enhancing multi-hop question answering.

# 6 ETHICAL CONSIDERATIONS

The knowledge homophily, where topologically proximate entities tend to share similar knowledgeability, increases the risk of knowledge-extraction attacks. An adversary can exploit homophily by crafting queries that target cohorts of related entities (or by combining entities across neighborhoods) to maximize coverage and force the model to reveal more information than intended. This raises privacy and copyright concerns because grouping semantically related entities makes it easier for the model to reconstruct or expose sensitive facts. To mitigate this risk, we propose defensive measures such as pre-filtering and sanitizing queries to remove suspicious or overly targeted prompts, rate-limiting repetitive/cohort queries, applying prompt-level heuristics that block requests for verbatim proprietary content, and deploying detector models or red-teaming workflows to identify adversarial extraction patterns. Combining these defenses with selective fine-grained access controls and differential privacy constraints can substantially reduce the attack surface introduced by knowledge homophily.

## REFERENCES

[1] Badr AlKhamissi, Millicent Li, Asli Celikyilmaz, Mona Diab, and Marjan Ghazvininejad. 2022. A review on language models as knowledge bases. *arXiv preprint arXiv:2204.06031* (2022).

[2] Weston A Bousfield and Charles Hill W Sedgewick. 1944. An analysis of sequences of restricted associative responses. *The Journal of General Psychology* (1944).

[3] Damai Dai, Li Dong, Yaru Hao, Zhifang Sui, Baobao Chang, and Furu Wei. 2021. Knowledge neurons in pretrained transformers. *arXiv preprint arXiv:2104.08696* (2021).

[4] Hady Elsahar, Pavlos Vougiouklis, Arslen Remaci, Christophe Gravier, Jonathon Hare, Frederique Laforest, and Elena Simperl. 2018. T-rex: A large scale alignment of natural language with knowledge base triples. In *Proceedings of the Eleventh International Conference on Language Resources and Evaluation (LREC 2018)*.

[5] Mikhail Galkin, Priyansh Trivedi, Gaurav Maheshwari, Ricardo Usbeck, and Jens Lehmann. 2020. Message passing for hyper-relational knowledge graphs. *arXiv preprint arXiv:2009.10847* (2020).

[6] Mor Geva, Roei Schuster, Jonathan Berant, and Omer Levy. 2020. Transformer feed-forward layers are key-value memories. *arXiv preprint arXiv:2012.14913* (2020).

[7] Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. 2020. Measuring massive multitask language understanding. *arXiv preprint arXiv:2009.03300* (2020).

[8] Mathieu Jacomy, Tommaso Venturini, Sebastien Heymann, and Mathieu Bastian. 2014. ForceAtlas2, a continuous graph layout algorithm for handy network visualization designed for the Gephi software. *PloS one* (2014).

[9] Shengbin Jia, Yang Xiang, Xiaojun Chen, and Kun Wang. 2019. Triple trustworthiness measurement for knowledge graph. In *The World Wide Web Conference*. 2865–2871.

[10] Saurav Kadavath, Tom Conerly, Amanda Askell, Tom Henighan, Dawn Drain, Ethan Perez, Nicholas Schiefer, Zac Hatfield-Dodds, Nova DasSarma, Eli Tran-Johnson, et al. 2022. Language models (mostly) know what they know. *arXiv preprint arXiv:2207.05221* (2022).

[11] Stephanie Lin, Jacob Hilton, and Owain Evans. 2021. Truthfulqa: Measuring how models mimic human falsehoods. *arXiv preprint arXiv:2109.07958* (2021).

[12] Bang Liu, Xinfeng Li, Jiayi Zhang, Jinlin Wang, Tanjin He, Sirui Hong, Hongzhang Liu, Shaokun Zhang, Kaitao Song, Kunlun Zhu, et al. 2025. Advances and challenges in foundation agents: From brain-inspired intelligence to evolutionary, collaborative, and safe systems. *arXiv preprint arXiv:2504.01990* (2025).

[13] Linhao Luo, Thuy-Trang Vu, Dinh Phung, and Gholamreza Haffari. 2023. Systematic assessment of factual knowledge in large language models. *arXiv preprint arXiv:2310.11638* (2023).

[14] Yao Ma, Xiaorui Liu, Neil Shah, and Jiliang Tang. 2021. Is homophily a necessity for graph neural networks? *International Conference on Learning Representations* (2021).

[15] Alex Mallen, Akari Asai, Victor Zhong, Rajarshi Das, Daniel Khashabi, and Hannaneh Hajishirzi. 2022. When not to trust language models: Investigating effectiveness of parametric and non-parametric memories. *arXiv preprint arXiv:2212.10511* (2022).

[16] Jeremy R Manning and Michael J Kahana. 2012. Interpreting semantic clustering effects in free recall. *Memory* (2012).

[17] Kevin Meng, David Bau, Alex Andonian, and Yonatan Belinkov. 2022. Locating and editing factual associations in gpt. *Advances in neural information processing systems* 35 (2022), 17359–17372.

[18] Xinyi Mou, Zejun Li, Hanjia Lyu, Jiebo Luo, and Zhongyu Wei. 2024. Unifying local and global knowledge: Empowering large language models as political experts with knowledge graphs. In *Proceedings of the ACM Web Conference 2024*. 2603–2614.

[19] Vishwas Mruthyunjaya, Pouya Pezeshkpour, Estevam Hruschka, and Nikita Bhutani. 2023. Rethinking language models as symbolic knowledge graphs. *arXiv preprint arXiv:2308.13676* (2023).

[20] Fabio Petroni, Tim Rocktäschel, Sebastian Riedel, Patrick Lewis, Anton Bakhtin, Yuxiang Wu, and Alexander Miller. 2019. Language Models as Knowledge Bases?. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*. 2463–2473.

[21] Pouya Pezeshkpour. 2023. Measuring and modifying factual knowledge in large language models. In *2023 International Conference on Machine Learning and Applications (ICMLA)*. IEEE, 831–838.

[22] Thorsten Rings, Timo Bröhl, and Klaus Lehnertz. 2022. Network structure from a characterization of interactions in complex systems. *Scientific Reports* 12, 1 (2022), 11742.

[23] Adam Roberts, Colin Raffel, and Noam Shazeer. 2020. How Much Knowledge Can You Pack Into the Parameters of a Language Model?. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*. 5418–5426.

[24] Tara Safavi and Danai Koutra. 2020. Codex: A comprehensive knowledge graph completion benchmark. *arXiv preprint arXiv:2009.07810* (2020).

[25] Yucheng Shi, Qiaoyu Tan, Xuansheng Wu, Shaochen Zhou, Kaixiong Zhou, and Ninghao Liu. 2024. Retrieval-enhanced knowledge editing in language models for multi-hop question answering. In *Proceedings of the 33rd ACM International Conference on Information and Knowledge Management*. 2056–2066.

[26] Nianwen Si, Hao Zhang, Heyu Chang, Wenlin Zhang, Dan Qu, and Weiqiang Zhang. 2023. Knowledge unlearning for llms: Tasks, methods, and challenges. *arXiv preprint arXiv:2311.15766* (2023).

[27] Zixing Song, Xiangli Yang, Zenglin Xu, and Irwin King. 2022. Graph-based semi-supervised learning: A comprehensive review. *IEEE Transactions on Neural Networks and Learning Systems* (2022).

[28] Olaf Sporns. 2012. The human connectome: a complex network. *Schizophrenia Research* (2012).

[29] Katherine Tian, Eric Mitchell, Allan Zhou, Archit Sharma, Rafael Rafailov, Huaxiu Yao, Chelsea Finn, and Christopher D Manning. 2023. Just ask for calibration: Strategies for eliciting calibrated confidence scores from language models fine-tuned with human feedback. *arXiv preprint arXiv:2305.14975* (2023).

[30] Denny Vrandečić and Markus Krötzsch. 2014. Wikidata: a free collaborative knowledgebase. *Commun. ACM* (2014).

[31] Yu Wang and Tyler Derr. 2021. Tree decomposed graph neural network. In *Proceedings of the 30th ACM international conference on information & knowledge management*. 2040–2049.

[32] Yu Wang, Nedim Lipka, Ryan A Rossi, Alexa Siu, Ruiyi Zhang, and Tyler Derr. 2024. Knowledge graph prompting for multi-document question answering. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 38. 19206–19214.

[33] Yu Wang, Yuying Zhao, Neil Shah, and Tyler Derr. 2022. Imbalanced graph classification via graph-of-graph neural networks. In *Proceedings of the 31st ACM international conference on information & knowledge management*. 2067–2076.

[34] Zhilin Yang, Peng Qi, Saizheng Zhang, Yoshua Bengio, William W Cohen, Ruslan Salakhutdinov, and Christopher D Manning. 2018. HotpotQA: A dataset for explainable multi-hop question answering. *arXiv preprint arXiv:1809.09600* (2018).

[35] Zijun Yao, Weijian Qi, Liangming Pan, Shulin Cao, Linmei Hu, Weichuan Liu, Lei Hou, and Juanzi Li. 2024. Seakr: Self-aware knowledge retrieval for adaptive retrieval augmented generation. *arXiv preprint arXiv:2406.19215* (2024).

[36] Zihan Zhang, Meng Fang, and Ling Chen. 2024. RetrievalQA: Assessing Adaptive Retrieval-Augmented Generation for Short-form Open-Domain Question Answering. In *Findings of the Association for Computational Linguistics ACL 2024*. 6963–6975.

[37] Danna Zheng, Mirella Lapata, and Jeff Z Pan. 2024. Large language models as reliable knowledge bases? *arXiv preprint arXiv:2407.13578* (2024).

[38] Shangshang Zheng, He Bai, Yizhe Zhang, Yi Su, Xiaochuan Niu, and Navdeep Jaitly. 2023. KGLens: Towards Efficient and Effective Knowledge Probing of Large Language Models with Knowledge Graphs. *arXiv preprint arXiv:2312.11539* (2023).

[39] Shuangjia Zheng, Jiahua Rao, Ying Song, Jixian Zhang, Xianglu Xiao, Evandro Fei Fang, Yuedong Yang, and Zhangming Niu. 2021. PharmKG: a dedicated knowledge graph benchmark for bomedical data mining. *Briefings in bioinformatics* (2021).

[40] Jiong Zhu, Ryan A Rossi, Anup Rao, Tung Mai, Nedim Lipka, Nesreen K Ahmed, and Danai Koutra. 2021. Graph neural networks with heterophily. In *Proceedings of the AAAI conference on artificial intelligence*, Vol. 35. 11168–11176.

**Figure 4: Homophily-guided Knowledge Injection and Retrieval: The process begins by training a GNN on a subset of entities with ground-truth knowledgeability scores (Blue Nodes) obtained by querying the base LLM. The trained GNN then infers the knowledgeability scores for all remaining entities (Green Nodes). Based on these predictions, triplets associated with entities estimated to have the lowest knowledge values are selected until the budget is met. Finally, the base LLM is fine-tuned on these less-known triplets to efficiently inject new knowledge, and its improved performance is measured on a held-out test set (Orange Nodes) in Figure 4(a). The estimated knowledgeability scores also guide retrieval, as illustrated in Figure 4(b).**

# A    APPENDIX

## A.1    Experimental Setting Visualization

*A.1.1    Homophily-guided Knowledge Injection.* Figure 4 (a) illustrates the pipeline of homophily-guided knowledge injection. In general, we leverage knowledge homophily to train a GNN estimator that predicts entity knowledgeability, and then use these estimated knowledgeability scores to identify less-known triplets for fine-tuning LLMs. The global procedure is as follows:

- **Step 1 - Fine-tuning Stage - Anchor Set Selection**: The process begins with a predefined triplet budget for fine-tuning LLMs. From this budget, 20% is allocated to anchor triplets, while the remaining 80% is reserved for knowledgeability estimation. Anchor triplets are constructed using an entity-centric sampling strategy: entities are randomly selected one by one, and all their associated triplets are added until the 20% quota is met. These anchor entities are used for training the GNN knowledge estimator and are shown as the blue "Known Knowledge Scores" nodes in Figure 4(a). Their ground-truth knowledgeability is obtained by querying the base LLM on their anchor triplets and aggregating the outcomes (see Section 3.1). With knowledge homophily, the GNN is then trained on these anchor entities to learn the relation between graph topology and knowledgeability.
- **Step 2 - Fine-tuning Stage - Estimating Knowledgeability of Remaining Set**: After training, the GNN estimator infers scores for all unlabeled entities (i.e., those with unknown knowledgeability). Entities are then ranked by their predicted knowledge value, $\mathcal{K}(v)$, where lower scores indicate a higher unknown level to the LLM. Finally, triplets linked to the least knowledgeable entities are selected as the remaining 80% of fine-tuning budget.
- **Step 3 - Evaluation Stage**: The selected triplets are combined with the initial anchor set to form the fine-tuning dataset. This

dataset, enriched with facts less known to the LLM, is used for fine-tuning. We evaluate the procedure in two ways. First, we measure selection quality, verifying whether the selected triplets are indeed unknown to the LLM. Second, we sample 2% of triplets (orange nodes) that are neither used in fine-tuning nor involve entities overlapped with fine-tuned triplets, to assess the generalization gain. As shown in Table 1, the fine-tuned LLM exhibits clear improvements in both selection quality and generalization performance. These results demonstrate the effectiveness of our homophily-aware knowledge injection in identifying knowledge deficiencies of LLMs and maximizing fine-tuning gains.

*A.1.2    Homophily-guided Knowledge Retrieval.* Figure 4(b) details the operational pipeline of our homophily-guided knowledge retrieval method in Section 4, designed to enhance the quality of the retrieved context to further improve multi-hop question answering.

The process begins with the creation of a multi-hop question set from 2-hop and 3-hop triplet paths. To ensure a fair evaluation, the entities that constitute these question paths are explicitly excluded from the GNN training data to prevent data leakage. From the remaining pool of entities, 40% are sampled to train the GNN regressor. This trained model infers the knowledge scores, $\mathcal{K}(v)$, for other entities, quantifying the awareness of LLMs of these other entities. For a given multi-hop question, the retrieval process commences with entity linking to anchor the query to a starting entity in the graph. From this point, a GNN-based Beam Search (G-BS) is employed to explore potential reasoning paths. The core of this method is its unique scoring function, $\mathcal{S} \times (1 - \alpha \cdot \mathcal{K}(u))$, which balances semantic relevance with a penalty based on the knowledgeability of the next entity ($\mathcal{K}(u)$). By penalizing the expansions toward well-known entities, the search prioritizes retrieving facts with higher information gain, thereby providing the LLM with specific context to answer the query.

**Table 3: Correlation results between entity knowledgeability scores under Full and Sparse settings across datasets.**

| Dataset | Full vs Sparse75% | | Full vs Sparse50% | |
|---|---|---|---|---|
| | Pearson | Spearman | Pearson | Spearman |
| **CoDEx-S** | 0.9577 | 0.9484 | 0.8490 | 0.8398 |
| **MVPKG** | 0.9376 | 0.9357 | 0.8839 | 0.8775 |
| **PharmKG** | 0.9444 | 0.9408 | 0.8773 | 0.8650 |
| **T-Rex** | 0.9177 | 0.9156 | 0.7957 | 0.7834 |
| **WD50K** | 0.9370 | 0.9278 | 0.8275 | 0.8080 |

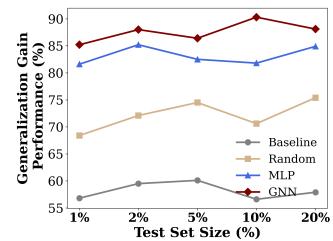## A.2 Additional Results

*A.2.1 Knowledgeability Score Robustness.* To validate that the calculated entity knowledgeability score ($\mathcal{K}(v)$) is robust to our triplet sampling algorithm, we conducted an experiment to assess its stability under data sparsity. For each knowledge graph, we created two sparsified versions: one retaining 75% of the original triplets (Sparse75%) and another retaining 50% (Sparse50%). We then recalculated the knowledgeability scores for all entities on these sparse graphs and computed the Pearson/Spearman correlations against the scores derived from the full, original graph. The results, presented in Table 3, demonstrate the stability of our metric.
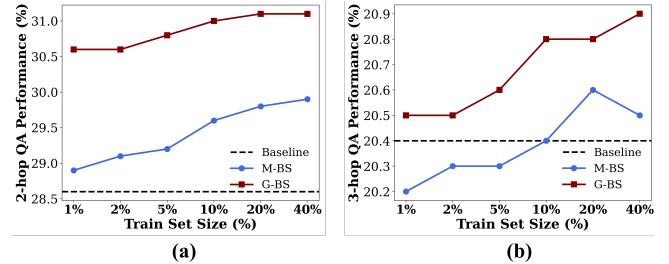
For the Sparse75% condition, the Pearson correlation consistently exceeded 0.91 across all datasets, indicating a very strong linear relationship. Even with half of the relational data removed in the Sparse50% condition, the scores maintained a strong correlation with the originals, with Pearson values generally above 0.80. The consistently high values for both Pearson and Spearman coefficients confirm that the entity knowledgeability score is not simply a byproduct of local graph density. These findings provide strong empirical evidence that our metric captures a stable property of the LLM's knowledge, which can be reliably estimated even from an incomplete set of relational facts.

*A.2.2 Sensitivity Analysis: Knowledge Injection.* To assess the robustness of our findings and ensure our conclusions are not dependent upon a specific test set size in the knowledge injection experiment, we conduct a sensitivity analysis. This analysis evaluates the performance of our fine-tuned models on the CoDEx-S dataset using the Mistral-7B model, mirroring the primary experiment in Section 4.2. We constructed five randomly sampled test sets, each representing a different proportion of the total dataset: 1%, 2%, 5%, 10%, and 20%. It was ensured that none of the entities used to train the GNN/MLP knowledgeability estimator appeared in any of the test sets. The results of this analysis are presented in the Figure 5. Our key findings of this analysis are:

- **Consistent Superiority:** The GNN-guided knowledge injection method consistently outperforms the MLP, Random, and Baseline approaches across all evaluation set sizes. This confirms the significant advantage of leveraging knowledge homophily.
- **Performance Stability:** Although minor fluctuations were observed, attributable to statistical variance in sampling, the performance of all methods remained relatively stable across evaluation set sizes. This finding suggests that the measured performance of the models is not merely a statistical artifact of a particular test set size.



**Figure 5: The knowledge injection performance of the fine-tuned Mistral models on the CoDEx-S dataset across varying test set sizes. The GNN-guided approach maintains a significant performance advantage over other methods.**



**Figure 6: The knowledge-aware retrieval performance across the varying training budgets of the underlying knowledge-ability estimator. For both 2-hop (left) and 3-hop (right) QA on the CoDEx-S dataset, the GNN-based search (G-BS) consistently outperforms the homophily-agnostic MLP-based search (M-BS) and the semantic baseline.**

*A.2.3 Sensitivity Analysis: Knowledge Retrieval.* To evaluate the robustness of our knowledge-aware beam search method (G-BS and M-BS), we conducted a sensitivity analysis on the amount of training data used for the knowledgeability estimators. The primary experiment in our paper utilizes GNN and MLP models trained on 40% of the available entities. This analysis investigates how performance on the multi-hop question-answering task varies when this training budget is reduced. For this sensitivity analysis we select the CoDEx-S dataset. We trained a series of GNN and MLP knowledgeability estimators on progressively larger subsets of entity data: 1%, 2%, 5%, 10%, 20%, and 40%. Each resulting estimator was then integrated into the G-BS and M-BS retrieval methods, respectively, and evaluated on the fixed set of 1,000 2-hop and 3-hop questions. The performance of the Semantic Beam Search baseline is independent of this training budget and remains constant. The results are presented in the Figure 6 for 2-hop (left) and 3-hop (right) QA performance. Our key findings of this analysis are:

- **Consistent Superiority of G-BS:** The GNN-based approach (G-BS) consistently outperforms both the homophily-agnostic MLP-based method (M-BS) and the baseline across all training set sizes and for both 2-hop and 3-hop questions. This confirms that the advantage of leveraging knowledge homophily is robust, even under data-scarce conditions.
- **Performance Scaling with Data:** The performance of both G-BS and M-BS improves as the training budget for the knowledge-ability estimator increases. This suggests that more accurately estimated knowledge scores lead to better path retrieval for QA.