

LARGE LANGUAGE MODELS COULD BE ROTE LEARNERS

Anonymous authors

Paper under double-blind review

ABSTRACT

Multiple-choice question (MCQ) benchmarks are widely used for evaluating Large Language Models (LLMs), yet their reliability is undermined by benchmark contamination. In this study, we reframe contamination as an inherent aspect of learning and seek to disentangle genuine capability acquisition from superficial memorization in LLM evaluation. First, by analyzing model performance under different memorization conditions, we uncover a counterintuitive trend: LLMs perform worse on memorized MCQs than on non-memorized ones, indicating the coexistence of two distinct learning phenomena, *i.e.*, rote memorization and genuine capability learning. To disentangle them, we propose **TrinEval**, a novel evaluation framework that reformulates MCQs into an alternative trinity format, reducing memorization while preserving knowledge assessment. Experiments validate TrinEval’s effectiveness in reformulation, and its evaluation reveals that common LLMs may memorize by rote 20.5% of knowledge points (in MMLU on average).

1 INTRODUCTION

The rapid advancement of Large Language Models (LLMs), driven primarily by large-scale pre-training on massive datasets, has endowed these models with remarkable proficiency across diverse tasks (Ouyang et al., 2022; OpenAI, 2024; Touvron et al., 2023). As LLMs continue to improve, evaluating their genuine capacities has emerged as one of the fundamental challenges, necessitating proper methodologies to ensure fairness and robustness (Ganguli et al., 2023; Liu et al., 2023b).

Among the developed methods, multiple-choice question (MCQ) benchmarks have become a standard approach for evaluation. Typically, LLMs are presented with a question and a fixed set of answer choices, requiring them to select the most appropriate option (see Fig. 1 for illustration). This format enables straightforward performance measurement through accuracy metrics and could cover a wide range of subjects. However, despite their widespread adoption, MCQ-based evaluation raises concerns about reliability due to benchmark contamination (Li & Flanigan, 2024; Kim et al., 2024), *i.e.*, test data unintentionally appears in training corpora and models may exploit memorized content rather than demonstrating genuine understanding, inflating their apparent capabilities. For instance, Zhou et al. (2023) discovers that smaller models with deliberate pre-exposure could outperform their larger counterparts, thereby contradicting widely accepted scaling laws.

To mitigate the issue, Zhou et al. (2023) advocates the removal of benchmark datasets from pre-training corpora. However, this strategy conflicts with the fundamental objective of large-scale pre-

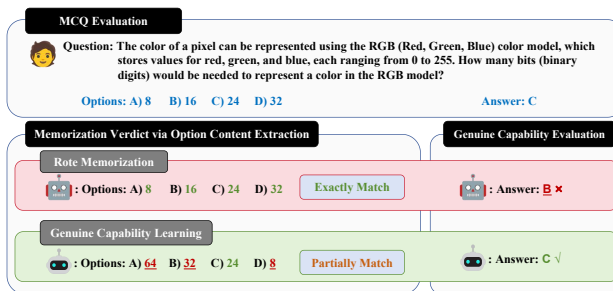


Figure 1: MCQ-based evaluation. We observe that LLMs tend to underperform on memorized MCQs.

054 training, which aims to maximize model performance by exposing LLMs to as much data as possi-
 055 ble. From a broader perspective, human learning also involves problem-solving through practicing
 056 on similar questions, *e.g.*, exam preparation. While rote memorization of specific questions and
 057 answers merely lead to short-term success, repeated practicing can also facilitate deeper conceptual
 058 understanding. Inspired, rather than viewing benchmark contamination as a flaw to be eradicated,
 059 which is a nearly impossible task at scale (Sainz et al., 2023; Bordt et al., 2024), we argue that it is an
 060 inherent aspect of learning and should be accounted for in evaluation. Therefore, this study shifts its
 061 focus to *evaluating LLMs in the presence of contamination, aiming to distinguish genuine capability*
 062 *gains from superficial memorization effects*. The explicit disentangling of these two learning effects
 063 remains largely unexplored in MCQ-based evaluation, yet we believe it marks a crucial step towards
 064 developing more rigorous and unbiased evaluation methodologies.

065 To investigate the effects of superficial memorization in LLM evaluation, we compare model per-
 066 formance under different memorization conditions. Inspired by membership inference attacks
 067 (MIA) (Carlini et al., 2022a; 2021), we define superficial memorization as an LLM’s ability to
 068 verbatim reproduce content, *e.g.*, exactly extracting option contents of MCQs in our case. Using
 069 this criterion, we partition the MMLU benchmark (Hendrycks et al., 2020)¹ into memorized and
 070 non-memorized subsets and evaluate three open-source LLMs² on both. Surprisingly, results reveal
 071 a consistent yet counterintuitive trend: LLMs perform worse on memorized MCQs than on those not
 072 (see Fig. 1 for illustration and Fig. 2 for results). This challenges the assumption that memorization
 073 improves model performance and suggests the coexistence of two distinct learning phenomena in
 074 LLMs: *rote memorization*, where models recall content verbatim without true understanding, and
 075 *genuine capability learning*, where they internalize underlying knowledge.

076 The preliminary investigation has several limitations. First, binary classification of MCQs as either
 077 memorized or non-memorized oversimplifies the nuances of memorization, potentially overlooking
 078 intermediate cases. Second, we rely on accuracy to measure performance, which is inherently unreli-
 079 able. Third, our analysis could not reveal the numerical dependency between rote memorization and
 080 capability learning. To address these challenges, we propose **TrinEval**, a novel evaluation frame-
 081 work designed to provide a more reliable measure of LLM performance by minimizing the influence
 082 of rote memorization. TrinEval employs a query-based probing (q-probing) mechanism Allen-Zhu
 083 & Li (2023) that reformulates MCQs into an alternative trinity format, *i.e.*, entity-attribute-context.
 This could prevent direct content recall while preserving knowledge assessment.

084 Through experiments, we demonstrate that TrinEval’s reformulation is knowledge-preserving, *i.e.*,
 085 maintaining testing problems’ inherent knowledge requirements without introducing extra cues, and
 086 could effectively reduce memorization. Combined with a continuous superficial memorization quan-
 087 tification metric, TrinEval reveals the in-robustness of LLMs’ capability learning, *e.g.*, with MMLU,
 088 tested open-sourced LLMs only mastered 19.6% of knowledge points while 20.5% are memorized
 089 by rote in the meanwhile, shedding light on the necessity for further optimization.

091 2 RELATED WORK

093 2.1 LLM EVALUATION ON MCQ BENCHMARKS

094
 095 The rapid advancement of LLMs has driven their expansion into diverse domains, necessitating
 096 robust and fair evaluation methodologies (Zheng et al., 2023b; Hu et al., 2025) and platforms (Con-
 097 tributors, 2023; Chiang et al., 2024). Among these, evaluating on MCQ benchmarks emerges as a
 098 widely adopted approach due to the ease of validation and standardized comparison across mod-
 099 els (Hendrycks et al., 2020; Wang et al., 2024; Zhong et al., 2023; Huang et al., 2024).

100 However, MCQ-based evaluations are not without limitations. Biases in LLM responses have been
 101 extensively studied (Dai et al., 2024), revealing issues such as social biases (Salewski et al., 2024;
 102 Liu et al., 2023a) and order sensitivity (Akter et al., 2023). To mitigate the latter, PriDe (Zheng
 103 et al., 2023a) estimates the option positional bias after option permutation. To examine mastery of
 104 knowledge, Zhao et al. (2023) applies a hypothesis testing method and checks rephrased-context
 105

106 ¹Selected for its popularity and documented data contamination in widely used LLMs (Sainz et al., 2023).

107 ²Llama2-7B (Touvron et al., 2023), Mistral-7B-v0.2 (Jiang et al., 2023) and Vicuna-v1.5-7B (Zheng et al., 2023b).

consistency for a given question. Benchmark contamination is arguably the most severe challenge for MCQ-based evaluations, which may result in misleadingly inflated performance (Zhou et al., 2023; Li & Flanagan, 2024). To address this, prior studies have explored data filtering, frequently-updated test sets (White et al., 2025), and data perturbation (Li et al., 2024).

In this paper, instead of attempting to eliminate contamination, we evaluate LLMs under its presence, aiming to distinguish genuine capability gains from superficial memorization effects. This marks a new perspective of LLM evaluation, revealing the extent to which models truly understand concepts rather than merely memorizing data.

2.2 LLM MEMORIZATION

Membership inference attacks (MIA) are commonly used to determine whether a specific sample was present in a model’s training data. Initially studied in smaller models, Carlini et al. (2022b) investigates deep learning memorization mechanisms by identifying and removing easily detectable memorized samples. In the context of LLMs, MIA has been employed to assess privacy risks, revealing that both open- and closed-source models can leak sensitive personal data when provided with related prompts (Kim et al., 2024).

Beyond privacy concerns, Carlini et al. (2022a) formally defines LLM memorization as a model’s ability to verbatim generate text sequences following a prefix prompt. Using this definition, several studies (Sainz et al., 2023; Bordt et al., 2024; Carlini et al., 2021) have examined mainstream LLMs, confirming widespread test data leakage across popular benchmarks. To quantify memorization strength, researchers (Shi et al., 2023; Zhang et al., 2024; Oren et al., 2023; Carlini et al., 2019) have further explored methods such as analyzing token probability distributions in generated outputs. However, while these studies extensively analyze LLM memorization, few explicitly investigate how memorization influences an LLM’s problem-solving ability. In contrast, our work focuses on their interplay, presenting a more rigorous approach to fair and reliable LLM evaluation.

3 METHODOLOGY

3.1 PRE-INVESTIGATION OF LLM CAPABILITY W.R.T. MEMORIZATION

Benchmark contamination often leads to inflated performance estimate. This phenomenon is commonly attributed to models memorizing specific questions and answers rather than demonstrating genuine problem-solving abilities. However, the extent to which and how memorization influences LLM performance remains unclear. To disentangle genuine capability acquisition from superficial memorization, we conduct a preliminary investigation into how LLMs perform under different memorization conditions. By examining model accuracy on memorized vs. non-memorized subsets, we aim to reveal the role of memorization in LLM evaluation and establish a foundation for more rigorous assessment methodologies.

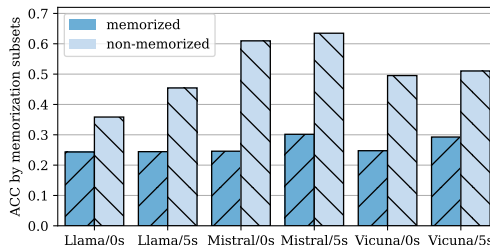


Figure 2: Model performance on memorized and non-memorized subsets of MMLU, where ‘0s’ and ‘5s’ stand for zero- and five-shot prompting, respectively.

In this paper, we mainly discussed the LLM Capability w.r.t. Memorization via MCQ problems as an example. The generalizability of our approach is discussed in the Appendix G. Formally, we define an MCQ as $x = \{x_Q, x_O, x_W\}$, where x_Q , x_O , and x_W refer to the question, options, and ground-truth answer, respectively. Following the memorization definition from Carlini et al. (2022a), we say an MCQ x is memorized by LLM G if G can extract/generate the content of options x_O exactly given question x_Q . In practice, given the commonly used next-token prediction pre-training, we incorporate meta-information (e.g., benchmark name) and 5-shot examples to elicit the memorized following content and recall memory (refer to Appendix E), and use greedy decoding (i.e., temperature fixed to 0) during extraction (Bordt et al., 2024; Sainz et al., 2023) (refer to Appendix A for the complete prompt). Using MMLU (Hendrycks et al., 2020) as the evaluation benchmark, we

162 divide the test set MCQs into memorized and non-memorized subsets, where the memorized subset
 163 consists of 909–982 questions (accounting for 6.5%–7.0% of the total 14,006) depending on the
 164 tested LLMs Llama2-7B, Mistral-7B-v0.2, and Vicuna-v1.5-7B. The detailed statistics of questions
 165 across subsets are given in Table 1 in Appendix, and we also observe that the majority of memorized
 166 questions are those relatively simple, *i.e.*, not in MMLU-PRO (Wang et al., 2024).

167 We then compute the accuracy (ACC) of tested LLMs by subsets as a proxy of model performance
 168 under different memorization conditions. The results of both zero- and five-shot prompting are re-
 169 ported in Fig. 2, from which we observe a consistent yet somehow counterintuitive trend: LLMs
 170 exhibit 47.2% lower accuracy on average on memorized MCQs compared to non-memorized ones,
 171 regardless of LLMs and prompting techniques. This finding challenges the commonly held as-
 172 sumption that memorization directly improves model performance. In addition, it also implies the
 173 coexistence of two distinct learning paradigms within LLMs, which we term rote memorization and
 174 genuine capability learning, respectively.

175 However, our pre-investigation has its limitations. The binary classification of memorization po-
 176 tentially overlooks more nuanced forms of learning. Additionally, using ACC as the performance
 177 metric does not truly capture model capacity. We address these two issues in the following subsec-
 178 tions, which then ensure a disentangle analysis of rote memorization and capability learning.

180 3.2 QUANTIFYING LLM MEMORIZATION

181 For quantifying the memorization of LLMs, prior research (Shi et al., 2023; Zhang et al., 2024)
 182 suggests that outlier tokens, which exhibit higher generation probabilities, are more likely to be
 183 found in memorized samples. Experiments on WIKIMIA (Shi et al., 2023) also finds that this
 184 method exhibits high AUC score on seen samples. Building on this idea, we develop a metric
 185 that utilizes the bottom $K\%$ of token probabilities within the generated sequence as a measure of
 186 memorization (we primarily set the value of K to 10 in this paper). Formally, the memorization
 187 score $F_m(\bar{x}, G)$ of LLM G on text sequence \bar{x} is computed as follows:

$$189 F_m(\bar{x}, G) = \frac{1}{|\mathcal{M}_K(\bar{x})|} \sum_{\bar{x}_i \in \mathcal{M}_K(\bar{x})} \log p_G(\bar{x}_i | \bar{x}_{1:i-1}), \quad (1)$$

191 where $p_G(\bar{x}_i | \bar{x}_{1:i-1})$ denotes the generation probability of token \bar{x}_i by G given its prefix subse-
 192 quence as context, and set $\mathcal{M}_K(\bar{x})$ includes the $K\%$ of tokens with the lowest probabilities. The
 193 higher F_m is, the more likely \bar{x} is memorized by the LLM, *i.e.*, the least memorized content could
 194 still be extracted with a high probability. Note that our objective is not to achieve high-precision
 195 quantification of memorization. Instead, we aim to statistically demonstrate a valid quantification of
 196 memorization across vast samples for comparative analysis.

198 3.3 MEASURING LLM CAPABILITY WITH TRINEVAL

200 We next present TrinEval, a novel evaluation framework designed to provide a more reliable measure
 201 of LLM performance by minimizing the influence of rote memorization.

202 To understand how LLMs store and manipulate knowledge, Allen-Zhu & Li (2023) created a fic-
 203 tional biography dataset that enumerates various attributes (*e.g.*, names, jobs, universities) and
 204 trained LLMs on this dataset. They employed a linear query-based probing method to uncover
 205 correlations between the entity token embeddings and the associated attributes, revealing that where
 206 LLMs encode knowledge, *e.g.*, under person names or sequence of the knowledge mention, is crucial
 207 for robust mastery of knowledge. This insight leads us to believe that entity tokens, which should
 208 ideally store related knowledge, are the target for evaluating an LLM’s genuine capability.

209 However, applying this method to real-world datasets, such as MMLU, presents challenges. Unlike
 210 controlled datasets with explicitly defined attributes, real-world data includes a far broader range
 211 of possible knowledge. As a result, we cannot enumerate all potential attributes and directly apply
 212 linear probing. To this end, we propose TrinEval, a verbal query probing method that reformulates
 213 MCQs around a knowledge-centric trinity: knowledge entity, attribute, and context. TrinEval is
 214 a pluggable augmentation on any MCQ benchmarks and could expose the genuine capability of
 215 LLMs by verifying whether they have correctly encoded knowledge. Generally, we instruct LLMs
 to extract the triplet according to the prompt and check and refine it if they are not qualified. The

reformulation is completed by a two-round reflection-based prompting method, with the detailed procedure (Alg. 1 in Appendix B) and related prompts available in Appendix B. Here, we explain the elements of the trinity and how to reformulate them.

Knowledge entity. We suppose that if an LLM has mastered some knowledge, the key information pertinent to the knowledge should be encoded within a few subject tokens, namely knowledge entity, to support efficient retrieval. By isolating these tokens, TrinEval ensures that only the essential information is considered.

Attribute. The attribute acts as a verbal probe to guide the model focusing on the specific feature or property of the knowledge entity being inquired. This mechanism allows TrinEval to isolate and assess the model’s understanding of the critical aspects of the questioning subject.

Context. In a certain portion of questions, the conditions or background context can significantly influence the solution approach. By explicitly including context in the evaluation process, TrinEval helps the model account for relevant situational details that might otherwise be overlooked, ensuring that the model’s answer is based on a comprehensive understanding of the problem.

By extracting the core and necessary question information in this trinity format, the reformulation by TrinEval is knowledge-preserving for the purposes of assessment. In the meanwhile, it completely destructs the original token sequence, effectively reducing the influence of memorization. We will empirically verify these properties through experiments.

Given an MCQ $x = (x_Q, x_O, x_W)$, it first queries a capable reformulation LLM to derive the knowledge entity x_E , attribute x_A , and context x_C from the original x . The LLM is instructed that the triplet should be sufficient for answering the question correctly, without including the answer option itself, ensuring the integrity of the evaluation. The same LLM then assesses whether the triplet contains all necessary information and no redundant information (typically, the rote-memorization), in the meanwhile, yields a rationale x_L as reflection (Shinn et al., 2024; Yao et al., 2022). If it does, the triplet is returned as the re-formulated question. Otherwise, the reformulation model refines the extraction, taking as input x_E, x_A, x_C , and x_L , and re-evaluates the updated triplet.

Finally, prompting with the extracted x_E, x_A, x_C , and options x_O , we inspect the generation probability of the next ground-truth answer token x_W (*i.e.*, A/B/C/D) as the measurement of capability:

$$F_c(x, G) = p_G(x_W | x_E, x_A, x_C, x_O). \quad (2)$$

As can be seen, F_c metric retains the necessary knowledge-centric information while discarding unnecessary biases, especially the rote memorization of LLMs, which leads to the quantification of genuine capability of LLMs.

4 EXPERIMENTS

In this section, we conduct extensive experiments to answer the following questions:

Q1. Does TrinEval effectively preserve knowledge to enable accurate knowledge assessment?

Q2. To what extent does TrinEval mitigate memorization effects during capability evaluation?

Q3. What insights does TrinEval provide into the distinction between rote memorization and genuine capability in LLMs?

4.1 EXPERIMENT SETUP

Models. We employ API-accessible commercial LLMs for question reformulation in TrinEval, specifically gpt-4o-2024-08-06 (GPT) (OpenAI, 2024) and qwen-max-2024-09-19 (Qwen) (Yang et al., 2024; Team, 2024). For model evaluation, we utilize open-source LLMs to obtain token-level logits. Our experiments focus on three widely adopted models: Llama2-7B (Llama) (Touvron et al., 2023), Mistral-7B-v0.2 (Mistral) (Jiang et al., 2023), and Vicuna-v1.5-7B (Vicuna) (Zheng et al., 2023b). These models are accessed from Hugging Face and implemented with the transformers library, allowing us to compute the log-probabilities of generated tokens for fine-grained analysis. All evaluations are conducted on a single NVIDIA A100 GPU with 40GB of memory. Default generation parameters are used, and greedy decoding is applied throughout to ensure reproducibility.

Benchmarks. We evaluate LLMs on the widely used MMLU benchmark (Hendrycks et al., 2020), in light of growing evidence of data contamination (Touvron et al., 2023) in many recent LLMs on this data. MMLU spans 57 subjects across STEM, humanities, social sciences, and others, providing a comprehensive assessment of model capabilities. We find duplicated MCQs across subjects on MMLU, and thus de-duplicate the dataset, resulting in a test set of 14,006 unique questions.

Evaluation. With commercial LLMs, we evaluate model performance by extracting predicted answers with regular expressions. For open-source models, we access the output probability of the first generated token (*i.e.*, corresponding to multiple-choice option IDs A/B/C/D) to compute quantitative performance metrics.

4.2 Q1. IS TRINEVAL KNOWLEDGE-PRESERVING?

We begin by verify whether TrinEval’s reformulation preserves essential knowledge, a prerequisite for reliable knowledge assessment. Specifically, our goal is to ensure that (1) the reformulation does not omit critical information that would cause originally correctly answered questions to be answered incorrectly, and (2) it does not introduce anomalous or unexpected content that could artificially inflate performance.

Following the complete TrinEval reformulation pipeline, we obtain 4,343 qualified MCQs with corresponding knowledge entities, attributes, and contexts using Qwen, and 4,645 MCQs with associated triplets using GPT. We then prompt both Qwen and GPT to answer their respective MCQs in both the original and reformulated triplet-based formats. Example MCQs and the questioning prompt templates in both formats are provided in Table 2 in the Appendix A; results are summarized in Fig. 3. Here, we primarily employ robust API-based LLMs to evaluate the knowledge retention capability of the TrinEval framework. This approach effectively mitigates potential interference from the models’ memorization effects, thereby ensuring accurate assessment of their intrinsic knowledge processing abilities and preventing erroneous conclusions that might arise from conflating memorization with genuine comprehension.

Among the qualified MCQs reformulated by Qwen and GPT, 4,133 and 4,409 are answered correctly in at least one of the two formats. Notably, the majority of these are correctly answered in both original and TrinEval formats, accounting for 90.3% with Qwen and 90.4% with GPT, indicating strong consistency in terms of problem-solving across formats. Only 121 questions for Qwen (2.9%) and 226 for GPT (5.1%) are answered correctly exclusively in the original format, suggesting the TrinEval reformulation does not omit essential information. Conversely, 278 questions for Qwen (6.7%) and 197 for GPT (4.5%) are answered correctly only in the reformulated format, which is comparable in general to the numbers for the original format alone. This further suggests that the reformulation does not introduce extra information that might artificially enhance model performance. Collectively, these findings provide strong evidence that TrinEval preserves the core knowledge necessary for answering, satisfying the requirements for reliable capability evaluation.

4.3 Q2. CAN TRINEVAL REDUCE MEMORIZATION?

In this subsection, we investigate whether the proposed TrinEval reformulation can suppress unnecessary memorization, thereby isolating and revealing the LLMs’ genuine capabilities under various circumstances. Following Bordt et al. (2024), we introduce dataset-specific cues into the questioning prompts and assess to what extent TrinEval mitigates memorization elicited by such prompts.

To evoke memorization, we embed dataset-related metadata into the input prompts, *i.e.*, the dataset name and in-context few-shot examples drawn from the same dataset (see Appendix E). For this set

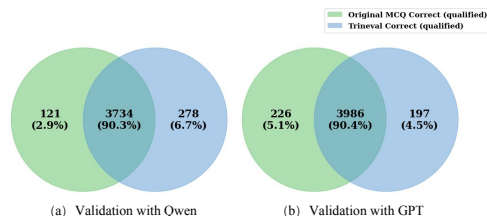


Figure 3: Knowledge-preserving validation of TrinEval reformulation. The green, blue, and overlapping regions represent the sets of MCQs correctly answered in the original format, TrinEval format, and both formats, respectively. Best viewed in color.

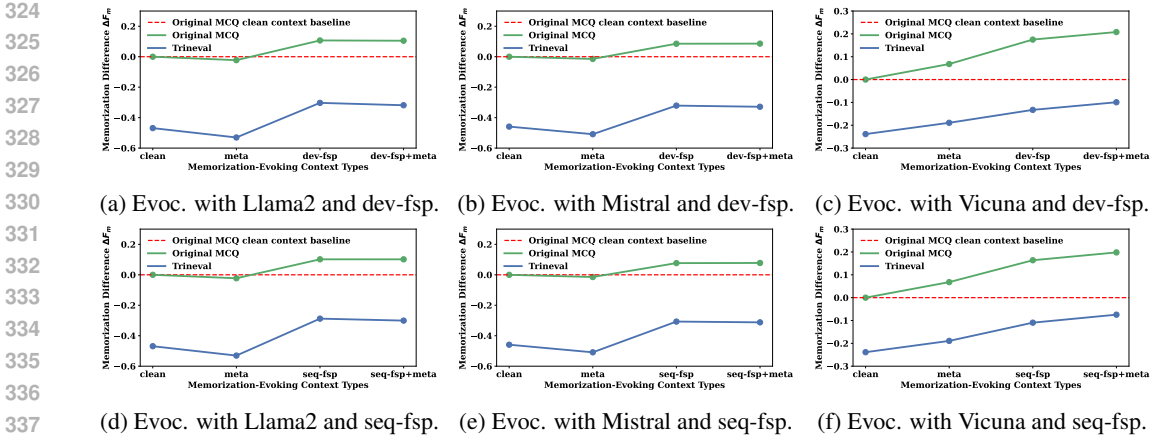


Figure 4: The results of memorization evocation (evoc.) under various dataset-related context, with green and blue curves referring to the memorization difference ΔF_m in the original and TrinEval formats, respectively. In the x-axis, ‘clean’, ‘meta’, ‘dev-fsp’, and ‘seq-fsp’ stand for without dataset-related context, with the name of the dataset, with few-shot prompt from the development set, and with few-shot prompt from the test set ahead of the testing question. These results of ΔF_m indicate the growing memorization effect given the increasing dataset-related information in general. However, the ΔF_m by TrinEval under the strongest memory evocation context remains consistently lower than the one in the original format, *e.g.*, the red dashed line.

of experiments, we focus on LLaMA, Mistral, and Vicuna, as their open-source implementations allow access to token-level output probabilities, enabling us to compute the memorization metric F_m . Since F_m lacks a defined absolute zero point indicating complete absence of memorization, we use the F_m value obtained from the original MCQ format without any additional context (*i.e.*, the vanilla MCQ) as a reference baseline. We then visualize the average change in F_m for each memorization-evoking method relative to this baseline.

Specifically, we hypothesize that unintentional data contamination may occur when LLMs are pre-trained on datasets scraped from public sources (*e.g.*, Hugging Face), where nearby examples from the same dataset may be concatenated during corpus construction. Following the pretraining mechanism of next-token prediction and prior findings from Carlini et al. (2022a), we posit that the inclusion of dataset-similar samples in the prompt may inadvertently trigger memorization. Accordingly, we design a sequence of memorization-evoking perturbations by progressively increasing the contextual cues: (1) providing only the dataset name, (2) adding few-shot examples from the same dataset (*i.e.*, samples in the development set or ahead examples in the test set of the same subject of MMLU), and (3) combining both types of context. For each MCQ, we compute the change in F_m relative to the reference baseline, capturing the degree to which memorization is elicited by specific context in original or TrinEval formats. Results are illustrated in Fig. 4.

Consistent with the observations by Bordt et al. (2024), the results show that F_m increases as more dataset-specific context is introduced, indicating stronger memorization effects. However, across all three tested LLMs, the ΔF_m curve for TrinEval remains consistently lower than that of the original MCQ format (*i.e.*, the red dashed line) regardless of various context used. Notably, even under the strongest memorization-evoking setting, TrinEval’s absolute F_m still stays below the vanilla baseline. These findings provide compelling evidence that TrinEval significantly mitigates the influence of memorization, effectively disentangling spurious recall from genuine model capability.

4.4 Q3. TRINEVAL’S FINDINGS ON LLM MEMORIZATION AND GENUINE CAPABILITY

In this subsection, we aim to explicitly examine the relationship between memorization and genuine problem-solving capability in LLMs, using the metrics F_m (memorization) and F_c (capability). Since commercial API-based models do not expose full vocabulary-level output probabilities, we focus on open-source LLMs for computing these metrics. After computing the F_m and F_c scores for all qualified MCQs in TrinEval format, we divide the values of F_m and F_c into five equal intervals,

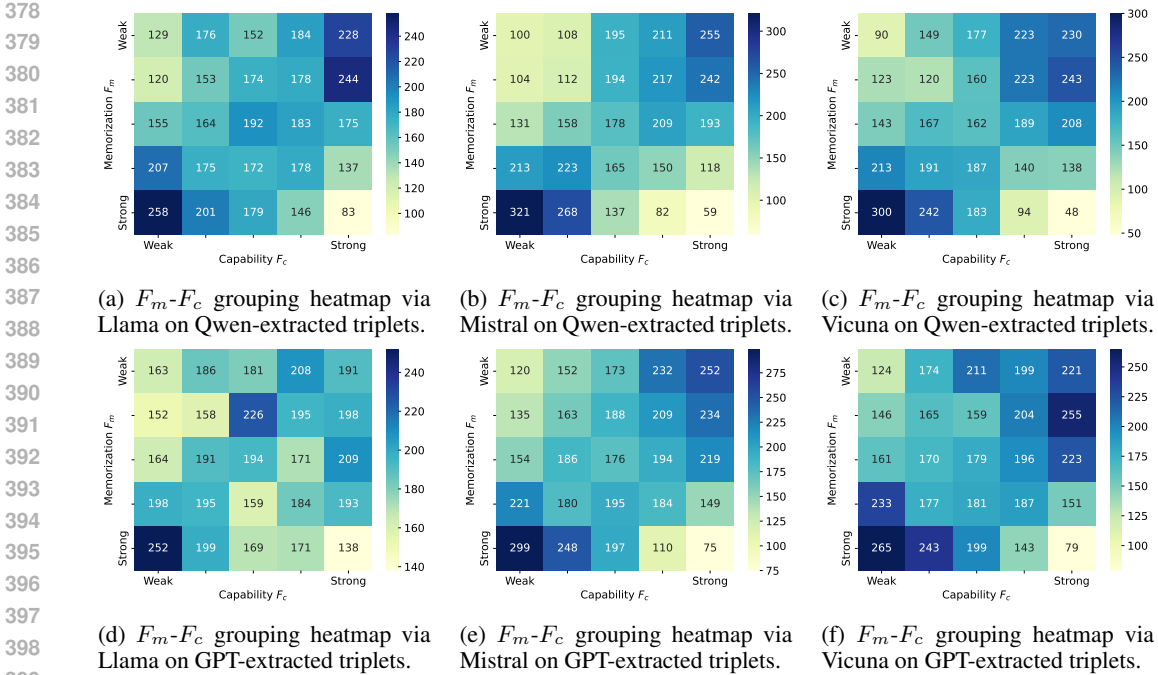


Figure 5: The distribution of MCQs based on memorization metric F_m vs. capability metric F_c . According to the values of F_m and F_c , we separate the MCQs into 25 groups and visualize the MCQ distribution from weak to strong with heatmaps.

respectively, resulting in 25 distinct groups of MCQ samples based on their joint distribution. We then use heatmaps to visualize the distribution of samples across these groups, thereby revealing the relationship between memorization and capability for each tested model.

As shown in Fig.5, the majority of MCQs cluster in the lower-left and upper-right corners of the heatmap. This pattern is evident and consistent across our tests. For instance, in the results of LLaMA using Qwen-extracted triplets, 38.57% of MCQs fall within the combined lower-left and upper-right 2×2 grid regions, yielding a Pearson correlation of -0.7755 (p-value < 0.05) for F_m vs. F_c of examples within the two square regions; expanding to the 3×3 regions increases coverage to 74.17%, with a stronger correlation of -0.8124 (p-value < 0.05). Similarly, for Mistral on Qwen triplets, 44.90% of MCQs lie in the two 2×2 corner regions with a correlation of -0.8722 (p-value < 0.05), and 80.82% in the two 3×3 corner regions with a correlation of -0.8794 (p-value < 0.05).

We take the MCQs within the lower-left 2×2 squares in Fig. 5 as the ones memorized by rote and compute the ratio of 20.5% by averaging the percentage of all evaluated LLMs. Similarly, 19.6% is computed as the ratio of MCQs within the upper-right 2×2 squares with genuine capability. Additional statistics results are reported in Table in Appendix C. These findings consistently indicate a negative correlation between memorization and capability: MCQs with lower memorization scores are more likely to reflect genuine problem-solving capability, whereas high memorization levels correspond to degraded performance, suggesting overfitting or shallow recall.

To interpret these results, we draw an analogy between LLM behavior and the human memory system, which comprises two key components: Short-Term Memory (STM) and Long-Term Memory (LTM) (Shiffrin, 2003). Neurobiological studies reveal that STM relies on transient synaptic protein synthesis with limited temporal persistence and functional scalability. In contrast, LTM is constructed through stabilized neuronal memory traces that constitute an enduring knowledge framework. This neural architecture not only supports STM operations as a cognitive substrate but also enables sophisticated information generalization across diverse contexts. This dichotomy aligns with recent observations in LLMs. As shown by Allen-Zhu & Li (2023) and Ovardia et al. (2023), models trained on corpora with diverse rephrasings exhibit stronger generalization than those trained solely on original formulations. When trained on a single fixed corpus format, LLMs, due to their

strong memorization capacity, tend to activate their STM system to memorize at the token level, capturing surface patterns rather than abstract knowledge. In this sense, **LLMs are potential rote learners**. Reformulations like TrinEval appear to facilitate the encoding of knowledge in a more principled, structured and reusable form, akin to LTM, thereby supporting more generalized and robust evaluation. Further experimental details and results are provided in Appendix F.

To further validate our hypothesis, we analyze the semantic proximity of MCQs within the qualified MMLU dataset by computing their embeddings and measuring the average distance to their closest 1% of samples. Our underlying assumption is that LLMs could better master a knowledge point if it is expressed in various formats in the training corpus. As a result, these knowledge ‘rephrases’ are also encoded densely via LLMs. By selecting only the top 1% nearest neighbors, we aim to retain the samples that are semantically-correlated in knowledge description. We then compare the mean embedding distances of MCQs located in the lower-left (strong memorization, weak capability) and upper-right (weak memorization, strong capability) 2×2 regions of the heatmap in Fig. 5. The corresponding results are presented in Fig. 6. Interestingly, we observe that the average embedding distance among the Genuine Capability Learning MCQs is much lower than (nearly half) that of the Rote Memorization MCQs. This suggests that memorized MCQs are more sparsely distributed in the embedding space, while non-memorized, capability-driven MCQs tend to form tighter semantic clusters. This observation also aligns well with the cognitive analogy of STM and LTM: rote memorization leads to fragmented, context-specific encoding, whereas genuine capability emerges from structured, reusable representations. The results of alternative thresholds are shown in Table in Appendix. Though it is well believed that memorization may lead to better but cheating performance of LLMs, we find that the more LLMs memorize, the worse they are at solving problems.

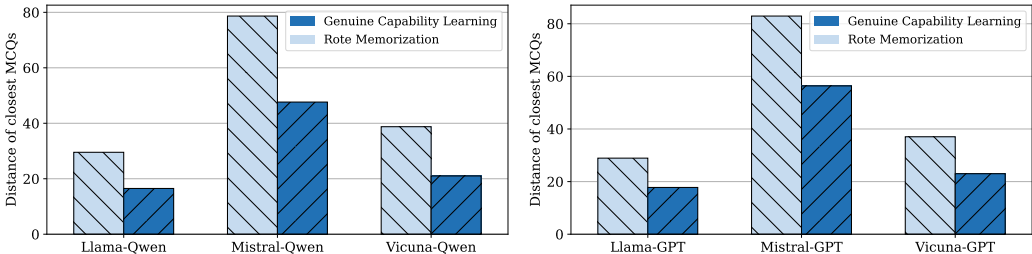


Figure 6: Averaged distance of each MCQs between the closest 1% MCQs’ embeddings. ‘Rote Memorization’ refers to MCQ within lower left 2×2 squares that typically exhibits high memorization metric F_m and low capability F_c while the ‘Genuine Capability Learning’ stands for MCQ lies within the upper right 2×2 squares with low F_m but high F_c . Further results are shown in Appendix F.

5 CONCLUSION

This study provided a novel perspective on benchmark contamination in LLM evaluation, reframing it as an inherent aspect of learning. This perspective led us to explore the relationship between memorization and genuine capability in LLMs. Through our empirical investigation, we observed a surprising result: LLMs performed worse on memorized MCQs compared to those not, suggesting that superficial memorization may undermine problem-solving ability rather than enhance it. This finding also implies the existence of two distinct learning paradigms in LLMs: rote memorization and genuine capability learning.

To disentangle them, we proposed TrinEval, an evaluation method reformulating MCQs into a knowledge-centric trinity and separating the influence of memorization from genuine knowledge application. Experiments validated both knowledge-preserving and memorization-reducing properties of TrinEval. Based on that, TrinEval reveals the in-robustness of LLMs’ knowledge learning, e.g., popular open-source LLMs memorize 20.5% of knowledge points by rote in MMLU. We also discussed the generalizability of our approach on other forms of questions (e.g., open-ended questions) in Appendix G. As such, we believe this work lays the groundwork for future studies on improving LLM knowledge robustness and more thorough evaluation.

REFERENCES

- 486
487
488 Syeda Nahida Akter, Zichun Yu, Aashiq Muhamed, Tianyue Ou, Alex Bäuerle, Ángel Alexander
489 Cabrera, Krish Dholakia, Chenyan Xiong, and Graham Neubig. An in-depth look at gemini’s
490 language abilities. *arXiv preprint arXiv:2312.11444*, 2023.
- 491 Zeyuan Allen-Zhu and Yuanzhi Li. Physics of language models: Part 3.1, knowledge storage and
492 extraction. *arXiv preprint arXiv:2309.14316*, 2023.
- 493 Sebastian Bordt, Harsha Nori, Vanessa Rodrigues, Besmira Nushi, and Rich Caruana. Elephants
494 never forget: Memorization and learning of tabular data in large language models. *arXiv preprint*
495 *arXiv:2404.06209*, 2024.
- 496
497 Nicholas Carlini, Chang Liu, Úlfar Erlingsson, Jernej Kos, and Dawn Song. The secret sharer:
498 Evaluating and testing unintended memorization in neural networks. In *28th USENIX security*
499 *symposium (USENIX security 19)*, pp. 267–284, 2019.
- 500 Nicholas Carlini, Florian Tramer, Eric Wallace, Matthew Jagielski, Ariel Herbert-Voss, Katherine
501 Lee, Adam Roberts, Tom Brown, Dawn Song, Úlfar Erlingsson, et al. Extracting training data
502 from large language models. In *30th USENIX Security Symposium (USENIX Security 21)*, pp.
503 2633–2650, 2021.
- 504 Nicholas Carlini, Daphne Ippolito, Matthew Jagielski, Katherine Lee, Florian Tramer, and
505 Chiyuan Zhang. Quantifying memorization across neural language models. *arXiv preprint*
506 *arXiv:2202.07646*, 2022a.
- 507
508 Nicholas Carlini, Matthew Jagielski, Chiyuan Zhang, Nicolas Papernot, Andreas Terzis, and Florian
509 Tramer. The privacy onion effect: Memorization is relative. *Advances in Neural Information*
510 *Processing Systems*, 35:13263–13276, 2022b.
- 511 Wei-Lin Chiang, Lianmin Zheng, Ying Sheng, Anastasios Nikolas Angelopoulos, Tianle Li,
512 Dacheng Li, Hao Zhang, Banghua Zhu, Michael Jordan, Joseph E Gonzalez, et al. Chatbot arena:
513 An open platform for evaluating llms by human preference. *arXiv preprint arXiv:2403.04132*,
514 2024.
- 515 OpenCompass Contributors. Opencompass: A universal evaluation platform for foundation models.
516 <https://github.com/open-compass/opencompass>, 2023.
- 517
518 Sunhao Dai, Chen Xu, Shicheng Xu, Liang Pang, Zhenhua Dong, and Jun Xu. Unifying bias and
519 unfairness in information retrieval: A survey of challenges and opportunities with large language
520 models. *arXiv preprint arXiv:2404.11457*, 2024.
- 521 Deep Ganguli, Nicholas Schiefer, Marina Favaro, and Jack Clark. Challenges in
522 evaluating AI systems, 2023. URL [https://www.anthropic.com/index/](https://www.anthropic.com/index/evaluating-ai-systems)
523 [evaluating-ai-systems](https://www.anthropic.com/index/evaluating-ai-systems).
- 524 Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and
525 Jacob Steinhardt. Measuring massive multitask language understanding. *arXiv preprint*
526 *arXiv:2009.03300*, 2020.
- 527
528 Renjun Hu, Yi Cheng, Libin Meng, Jiaxin Xia, Yi Zong, Xing Shi, and Wei Lin. Training an llm-
529 as-a-judge model: Pipeline, insights, and practical lessons. 2025.
- 530 Yuzhen Huang, Yuzhuo Bai, Zhihao Zhu, Junlei Zhang, Jinghan Zhang, Tangjun Su, Junteng Liu,
531 Chuancheng Lv, Yikai Zhang, Yao Fu, et al. C-eval: A multi-level multi-discipline chinese eval-
532 uation suite for foundation models. *Advances in Neural Information Processing Systems*, 36,
533 2024.
- 534 Albert Q Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot,
535 Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, et al.
536 Mistral 7b. *arXiv preprint arXiv:2310.06825*, 2023.
- 537
538 Siwon Kim, Sangdoon Yun, Hwaran Lee, Martin Gubri, Sungroh Yoon, and Seong Joon Oh. Propile:
539 Probing privacy leakage in large language models. *Advances in Neural Information Processing*
Systems, 36, 2024.

- 540 Changmao Li and Jeffrey Flanigan. Task contamination: Language models may not be few-shot any-
541 more. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pp. 18471–
542 18480, 2024.
- 543 Jiatong Li, Renjun Hu, Kunzhe Huang, Yan Zhuang, Qi Liu, Mengxiao Zhu, Xing Shi, and Wei
544 Lin. Perteval: Unveiling real knowledge capacity of LLMs with knowledge-invariant perturba-
545 tions. In *The Thirty-eight Conference on Neural Information Processing Systems Datasets and*
546 *Benchmarks Track*, 2024. URL <https://openreview.net/forum?id=HB5q6pC5eb>.
- 547 Yanchen Liu, Srishti Gautam, Jiaqi Ma, and Himabindu Lakkaraju. Investigating the fairness of
548 large language models for predictions on tabular data. *arXiv preprint arXiv:2310.14607*, 2023a.
- 549 Yang Liu, Yuanshun Yao, Jean-Francois Ton, Xiaoying Zhang, Ruocheng Guo Hao Cheng, Yegor
550 Klochkov, Muhammad Faaiz Taufiq, and Hang Li. Trustworthy llms: A survey and guideline for
551 evaluating large language models’ alignment. *arXiv preprint arXiv:2308.05374*, 2023b.
- 552 OpenAI. Hello gpt-4o, May 2024. URL <https://openai.com/index/hello-gpt-4o/>.
- 553 Yonatan Oren, Nicole Meister, Niladri Chatterji, Faisal Ladhak, and Tatsunori B Hashimoto. Proving
554 test set contamination in black box language models. *arXiv preprint arXiv:2310.17623*, 2023.
- 555 Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong
556 Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to fol-
557 low instructions with human feedback. *Advances in neural information processing systems*, 35:
558 27730–27744, 2022.
- 559 Oded Ovadia, Menachem Brief, Moshik Mishaeli, and Oren Elisha. Fine-tuning or retrieval? com-
560 paring knowledge injection in llms. *arXiv preprint arXiv:2312.05934*, 2023.
- 561 Oscar Sainz, Jon Ander Campos, Iker García-Ferrero, Julen Etxaniz, Oier Lopez de Lacalle, and
562 Eneko Agirre. Nlp evaluation in trouble: On the need to measure llm data contamination for each
563 benchmark. *arXiv preprint arXiv:2310.18018*, 2023.
- 564 Leonard Salewski, Stephan Alaniz, Isabel Rio-Torto, Eric Schulz, and Zeynep Akata. In-context im-
565 personation reveals large language models’ strengths and biases. *Advances in Neural Information*
566 *Processing Systems*, 36, 2024.
- 567 Weijia Shi, Anirudh Ajith, Mengzhou Xia, Yangsibo Huang, Daogao Liu, Terra Blevins, Danqi
568 Chen, and Luke Zettlemoyer. Detecting pretraining data from large language models. *arXiv*
569 *preprint arXiv:2310.16789*, 2023.
- 570 RM Shiffrin. Chapter: Human memory: A proposed system and its control processes. *Spence, KW;*
571 *Spence, JT. The psychology of learning and motivation. New York*, 2003.
- 572 Noah Shinn, Federico Cassano, Ashwin Gopinath, Karthik Narasimhan, and Shunyu Yao. Reflexion:
573 Language agents with verbal reinforcement learning. *Advances in Neural Information Processing*
574 *Systems*, 36, 2024.
- 575 Qwen Team. Qwen2.5: A party of foundation models, September 2024. URL [https://qwenlm.](https://qwenlm.github.io/blog/qwen2.5/)
576 [github.io/blog/qwen2.5/](https://qwenlm.github.io/blog/qwen2.5/).
- 577 Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Niko-
578 lay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open founda-
579 tion and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023.
- 580 Yubo Wang, Xueguang Ma, Ge Zhang, Yuansheng Ni, Abhranil Chandra, Shiguang Guo, Weiming
581 Ren, Aaran Arulraj, Xuan He, Ziyang Jiang, et al. Mmlu-pro: A more robust and challenging
582 multi-task language understanding benchmark. *arXiv preprint arXiv:2406.01574*, 2024.
- 583 Colin White, Samuel Dooley, Manley Roberts, Arka Pal, Benjamin Feuer, Siddhartha Jain, Ravid
584 Shwartz-Ziv, Neel Jain, Khalid Saifullah, Sreemanti Dey, Shubh-Agrawal, Sandeep Singh
585 Sandha, Siddhartha Venkat Naidu, Chinmay Hegde, Yann LeCun, Tom Goldstein, Willie
586 Neiswanger, and Micah Goldblum. Livebench: A challenging, contamination-free LLM bench-
587 mark. In *The Thirteenth International Conference on Learning Representations*, 2025. URL
588 <https://openreview.net/forum?id=sKYHBTaxVa>.

- 594 An Yang, Baosong Yang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Zhou, Chengpeng Li,
595 Chengyuan Li, Dayiheng Liu, Fei Huang, et al. Qwen2 technical report. *arXiv preprint*
596 *arXiv:2407.10671*, 2024.
- 597 Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik Narasimhan, and Yuan Cao.
598 React: Synergizing reasoning and acting in language models. *arXiv preprint arXiv:2210.03629*,
599 2022.
- 600 Jingyang Zhang, Jingwei Sun, Eric Yeats, Yang Ouyang, Martin Kuo, Jianyi Zhang, Hao Frank
601 Yang, and Hai Li. Min-k%++: Improved baseline for detecting pre-training data from large
602 language models. *arXiv preprint arXiv:2404.02936*, 2024.
- 603 Yukun Zhao, Lingyong Yan, Weiwei Sun, Guoliang Xing, Chong Meng, Shuaiqiang Wang, Zhicong
604 Cheng, Zhaochun Ren, and Dawei Yin. Knowing what llms do not know: A simple yet effective
605 self-detection method. *arXiv preprint arXiv:2310.17918*, 2023.
- 606 Chujie Zheng, Hao Zhou, Fandong Meng, Jie Zhou, and Minlie Huang. Large language models
607 are not robust multiple choice selectors. In *The Twelfth International Conference on Learning*
608 *Representations*, 2023a.
- 609 Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang,
610 Zi Lin, Zhuohan Li, Dacheng Li, Eric Xing, et al. Judging llm-as-a-judge with mt-bench and
611 chatbot arena. *Advances in Neural Information Processing Systems*, 36:46595–46623, 2023b.
- 612 Wanjun Zhong, Ruixiang Cui, Yiduo Guo, Yaobo Liang, Shuai Lu, Yanlin Wang, Amin Saied,
613 Weizhu Chen, and Nan Duan. Agieval: A human-centric benchmark for evaluating foundation
614 models. *arXiv preprint arXiv:2304.06364*, 2023.
- 615 Kun Zhou, Yutao Zhu, Zhipeng Chen, Wentong Chen, Wayne Xin Zhao, Xu Chen, Yankai Lin,
616 Ji-Rong Wen, and Jiawei Han. Don’t make your llm an evaluation benchmark cheater. *arXiv*
617 *preprint arXiv:2311.01964*, 2023.
- 618 Kaijie Zhu, Jiaao Chen, Jindong Wang, Neil Zhenqiang Gong, Diyi Yang, and Xing Xie. Dyval:
619 Dynamic evaluation of large language models for reasoning tasks. In *The Twelfth International*
620 *Conference on Learning Representations*, 2023.
- 621 Kaijie Zhu, Jindong Wang, Qinlin Zhao, Ruochen Xu, and Xing Xie. Dyval 2: Dynamic evaluation
622 of large language models by meta probing agents. *arXiv preprint arXiv:2402.14865*, 2024.

623 A DETAILS OF THE EXTRACTING PROMPTS AND THE EXTRACTED 624 (MEMORIZED) MCQS

625 In this section, we introduce the details of the processed dataset and the prompts for extraction.
626 MCQs from some subjects contain similar or identical options³. With the provided 5-shot prompt,
627 options of MCQs from these subjects can be easily extracted, leading to a high False-positive ratio.
628 In order to avoid the influence of the few-shot prompt on the option extraction, we eliminate MCQs
629 in which any of the options have appeared twice in the dataset. After deduplication, we obtain
630 14,006 MCQs for evaluation. The extraction prompt and the detailed statistics are shown in the
631 following text and Tab. 1.

632 B DETAILS OF TRINEVAL

633 In this section, we introduce the details of the proposed TrinEval. The pseudo-code is shown in
634 the Alg. 1. The prompts used are also shown below. Note that the potential data leakage is often
635 caused by the data crawled on the Huggingface dataset site. Thus, we also provide the original
636 text of MCQs in the format on the Huggingface dataset site to mimic the data contamination with
637 in-context learning.

638 ³E.g., the options of MCQs in the subject, moral_scenarios, are all identical (‘Wrong, Wrong’, ‘Wrong, Not
639 wrong’, ‘Not wrong, Wrong’ and ‘Not wrong, Not wrong’).

Prompt template for extraction:	Prompt template for pre-investigation on LLM Memorization w.r.t. Capability:
<p><i>You are an expert of multiple choice questions of MMLU dataset. The following are multiple-choice questions (with answers) about [subject].</i></p> <p><i>[examples]</i></p> <p><i>[question]</i></p> <p><i>Options:</i></p> <p><i>A.</i></p>	<p><i>You are an expert of multiple choice questions of MMLU dataset. The following are multiple choice questions (with answers) about [subject].</i></p> <p><i>[examples]</i></p> <p><i>[question]</i></p> <p><i>Options:</i></p> <p><i>A. [content for option A]</i></p> <p><i>B. [content for option B]</i></p> <p><i>C. [content for option C]</i></p> <p><i>D. [content for option D]</i></p> <p><i>Answer:</i></p>

Model	Subset	Simple	Pro	MMLU
Llama	memorized	912	70	982
	non-mem.	6,548	6,476	13,024
	all	7,460	6,546	14,006
Mistral	memorized	879	36	915
	non-mem.	6,581	6,510	13,091
	all	7,460	6,546	14,006
Vicuna	memorized	893	16	909
	non-mem.	6,567	6,530	13,097
	all	7,460	6,546	14,006

Table 1: Statistics of memorized and non-memorized questions by Llama2-7B, Mistral-7B-v0.2, and Vicuna-v1.5-7B in MMLU.

C DETAILED RESULTS OF MEMORIZATION V.S. CAPABILITY

In this section, we exhibit the detailed results of the Q3. What does TrinEval reveal about the memorization v.s. the capability of LLMs. We reveal the ratio of MCQs within the upper right and lower left 2×2 and 3×3 squares as well as the Pearson correlations between the F_m and F_c of these MCQs. Our analysis reveals a tendency towards a negative correlation between the capabilities and memorization of LLMs shown in the Tab. 3. All person correlation values are computed with p-values at the 0.05 level.

Further, inspired by the Precision-Recall Curve, we take each unique F_m of the qualified MCQs as the threshold to separate them as the Memorized and Capable MCQs. For each separation, we compute the probability of whether the F_c of a randomly selected Capable MCQ exceeds the F_c of a randomly selected Memorized MCQ and plot them as the blue curve. We also compute the T-test p-value between the F_c s of the Memorized MCQs and Capable MCQs as the green curve. The results are shown in Fig. 7. For the second row, we filter out the MCQs within the upper left and lower right 2×2 squares. From the figure, we observe that over a relatively long segment in the middle of the x-axis threshold range, the probability remains at a comparatively high value, while the p-value stays below 0.05. From this, we can conclude that F_m can distinguish between MCQs with high F_c and those with low F_c with a negative correlation at a high confidence level. This further supports that LLMs are potential rote learners, the more the LLMs memorize, the more poorly they perform.

Algorithm 1 MCQ reformulation by TrinEval

Input: Question x_Q , options x_O , and answer x_W of an MCQ.
Output: Reformulated question x_Q^R .
 1: Preliminarily extract knowledge entity x_E , attribute x_A , and context x_C based on x_Q, x_O and x_W ;
 2: Initialize $X_Q^R = x_E, x_A, x_C$;
 3: Validate the adequacy and necessity of the x_Q^R and give reasons x_L ;
 4: **if** x_Q^R matches the requirement **then**
 5: Return x_Q^R ;
 6: **else**
 7: Re-extract $x'_E, x'_A,$ and x'_C by reflecting with x_E, x_A, x_C and x_L ;
 8: Update $x_Q^R = x'_E, x'_A, x'_C$;
 9: Validate the adequacy and necessity of the x_Q^R and give reasons x_L ;
 10: **if** x_Q^R matches the requirement **then**
 11: Return x_Q^R ;
 12: **else**
 13: Discard the MCQ, return *None*;
 14: **end if**
 15: **end if**

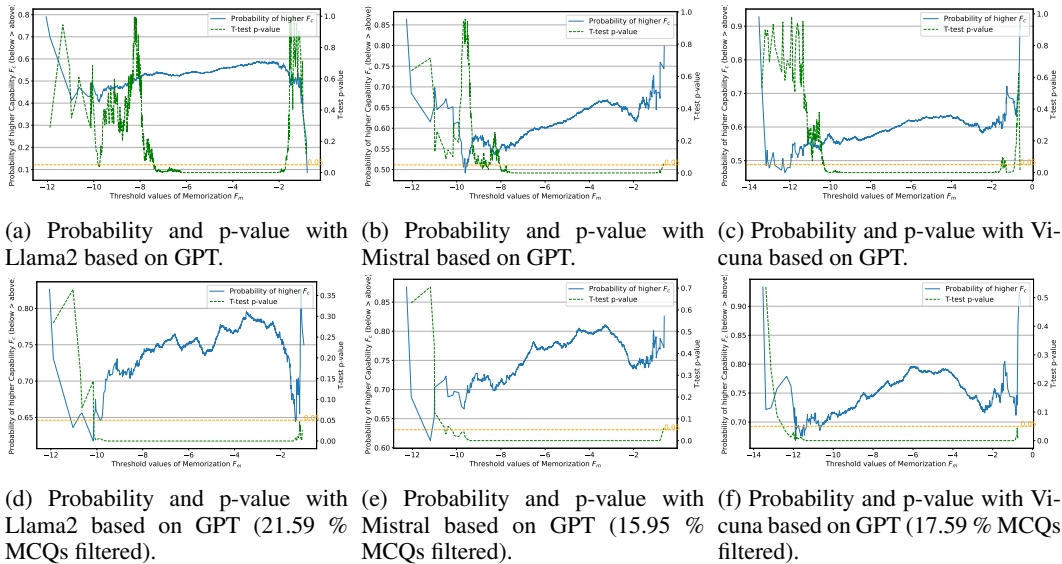


Figure 7: The over-performing probability curve and p-value curve with different F_m thresholds. In this figure, we take each unique F_m as the threshold to separate the qualified MCQs into the Memorized and non-Memorized MCQs. We compute the probability of a randomly selected non-Memorized MCQ’s F_c exceeding a randomly selected Memorized MCQ’s F_c under each threshold as the blue curve, and the green curve is the p-value of the T-test between the F_c s of the non-Memorized MCQs and the Memorized MCQs.

D HUMAN ANNOTATION

As there is potential risk for the LLM-based Knowledge Preserving evaluation in TrinEval procedure (line 4 and line 10 in 1) that the API-based LLMs might still be able to answer the MCQs without sufficient knowledge since this is still prompting LLMs who have been trained on these datasets and “know” the original content, human annotation is also applied. The annotation of each MCQ encompasses three subtasks: (1) answering the question in the re-organized form, (2) answering the question in the original form, and (3) verifying if the reformulation is Knowledge Preserving or not.

For efficient annotation, we implemented a stratified sampling procedure by selecting one MCQ per subject from all 56 MMLU subjects (as a temporary compromise for limited time, which will be expanded later) under both Qwen and GPT reorganization paradigms. This yielded 112 representative questions (2 systems \times 56 subjects) for evaluation. Three human annotators independently performed dual-form assessments through: (1) Direct question answering with the reformed format first and the original format; (2) Knowledge Preservation (K.P.) scoring across two dimensions: i. Knowledge adequacy (sufficiency for accurate response), ii. removal of redundant content using a 5-point scale (1=unsatisfactory; 2=major information are missed or unnecessary information is incorporated, but part is still acceptable, 3=need to take some time to understand, but can still solve the MCQ, 4=an element properly belonging to one triplet component appears in another, but does not impact MCQ solving; 5=optimal). We use a continuous rather than binary metric to mitigate the cognitive difference of the threshold between the annotators. Inter-rater reliability was ensured through consensus-building discussions prior to formal annotation. Final scores were aggregated using mean values to further mitigate individual annotator bias. The results are shown in Tab. 4 below.

Notably, our analysis reveals that over 95% of correctly answered MCQs maintained consistency across both original and paraphrased formats. Furthermore, human annotators rated our paraphrased questions mean K.P. scores exceeding 4.0 (on a 5-point scale), which means that the reformulated MCQs only somehow influence the readability of humans but do not impact the solvability of the original format. This provides empirical validation that our proposed TrinEval methodology effectively preserves necessary knowledge elements from original MCQ formulations, while the influence of the LLM memorization during the evaluation is rather limited.

Experimental results under human annotation also reveal that Qwen underperforms GPT in key metrics, particularly in processing long-context texts where it occasionally omits background information (evidenced by excessive "N/A" assignments in the Context fields). This capability gap is further reflected in MCQ annotations: there are only MCQs that are merely correctly answered with the original format except for the correct MCQs with both formats.

E ELABORATION ON DATASET-RELATED INFORMATION

We suppose that unintentional data contamination arises from crawling dataset pages (e.g., Hugging Face) during the compilation of LLM pretraining datasets. When researchers are organizing the pretraining corpus, one or more neighboring original data samples would be truncated and concatenated sequentially into a pretraining sample. Thus, according to Carlini’s theory Carlini et al. (2022a) and the next-token-prediction pretraining, we believe that offering samples within the same dataset would affect memorization evocation.

Besides, the previous study Bordt et al. (2024) also applies a similar method, the “Header Completion Test”, for tabular data memorization detection. By offering the heading rows within a CSV file, they also find that providing the preceding data samples can help to detect the memorized dataset by LLMs, which also practically proves that offering samples within the same dataset would affect the memorization evocation.

Following this, similar to the dataset name, we take the samples within the same dataset as the few-shot prompt in order to find out if the TrinEval re-organization method can avoid such memorization evocation phenomenon. Still, in Fig. 4, we can see that the blue curves remain below the green ones for the “def-fsp” setting, which proves that TrinEval can restrain the memorization evocation.

F EMBEDDING DISTANCE OF MEMORIZED AND NON-MEMORIZED MCQS

As there are 57 different subjects within the MMLU dataset, we believe unrelated sequences would lead to increased embedding distance even though they are among the mastered knowledge points. Here, we try to filter out the unrelated samples for each sample and thus filter the closest samples at the $1e - 2$ level in order to make sure there are not too many unrelated samples incorporated.

To highlight this result more prominently, we employed a relatively stringent data filtering strategy in the paper and made it 1% of the closest samples in Section 4.4. In the following version, we will add this clarification part in the paper. Here in order to provide a more robust result, we also

810 present results obtained under more lenient data filtering criteria, such as thresholds of 3% and 5%.
811 The results are shown in Tab. 5 (RM stands for rote memorization, and GCL stands for genuine
812 capability learning).

813 We can see that, as we said above, the more samples we incorporated, the higher the average distance
814 of the closest embeddings grows. Still, though we increase the filtering threshold and the distance
815 gap between the RM and the GCL is narrowing, we can still find that the distance between the
816 closest rote memorization MCQ embeddings is more than the distance between the closest genuine
817 capability learning MCQ embeddings. This proves that the reported results are robust and solid.
818

819 G LIMITATIONS

820
821 Our limitations are mainly three points. First, though we mainly use the MCQs as the tested bench-
822 marks, We believe that rote memorization and genuine capability learning are among the most essen-
823 tial for understanding LLM learning. In the current study, we choose MCQs as the testbed due to its
824 wide adoption in LLM evaluation and easily verifiable response correctness. This has led to several
825 innovations, such as the attempt to evaluate LLMs in the presence of contamination, the interintu-
826 itive trend of model performance under different memorization conditions, and the quantification of
827 rote memorization and genuine capability learning.

828 On the other hand, our disentangling point of view and the concrete methodology should also apply
829 to problems in forms other than MCQs, e.g., open-ended QA problems. Note that it could be possible
830 to recruit human experts to write responses to questions from these classic evaluation benchmarks,
831 says TruthfulQA, SuperGELU, Arena-Hard, and AlpacaEval 2.0. Using these responses during SFT
832 could also result in inflated performance. It would be very interesting to compare the performance
833 before and after question reformulation that reduces rote memorization to intentionally injected
834 responses. We will discuss such potential applications in the next version. Thank you again for your
835 comment.

836 Second, though our proposed TrinEval retrains the problem-solving ability of the LLMs and obtains
837 stronger robustness, it is not a dynamical re-organizing method that can still be leaked and pre-
838 experienced during training. On the one hand, we appeal to the LLM developers not to use this
839 re-organizing method as part of the training corpus. On the other hand, future works will be focused
840 on developing dynamic evaluation method (Zhu et al., 2023; 2024).

841 Finally, we did not give a clear exploration on how and why the more LLMs memorize, the less the
842 capability of the LLMs obtains. In future work, we will also look into the mechanism of the training
843 and structure of LLMs for a thorough study of the phenomenon.
844
845
846
847
848
849
850
851
852
853
854
855
856
857
858
859
860
861
862
863

864
865
866
867
868
869
870
871
872
873
874
875
876
877
878
879
880
881
882
883
884
885
886
887
888
889
890
891
892
893
894
895
896
897
898
899
900
901
902
903
904
905
906
907
908
909
910
911
912
913
914
915
916
917

Prompt template for triplet extraction:

You are an expert of Knowledge Keyword extraction. Analyze and summarize the Question based on the given Fact corpus and extract the Knowledge Keyword, the Attribute and the Context (if necessary) within the Question.

Given a Fact corpus, a Question about the Fact corpus, and the Answer to the Question, analyze the Question corpus as well as the given Answer. Applying the provided steps, extract the Knowledge Keyword, the Attribute of the Knowledge Keyword and the necessary Context to obtain the key information of the Question, ensuring they are sufficient for answering the given Question and obtaining the given Answer.

Steps

- 1. ****Review the Fact corpus:**** Read through the entire Fact corpus to understand the context.*
- 2. ****Identify the Question:**** Focus on the given Question to capture which part of the Fact corpus it is asking about.*
- 3. ****Understand the Answer to the Question:**** Compare the given Answer and the identified questioned part within the Fact corpus and understand why this answer was chosen.*
- 4. ****Write Step-by-Step Reasoning:*****
 - Identify the asked Knowledge Keyword in the Question that is the subject of the most information in the Fact corpus and the asked Question is about the information among.*
 - Determine the asked Attribute of the Knowledge Keyword in the Question, which can be used to infer the given Answer.*
 - Review the identified Knowledge Keyword and Attribute to confirm that only these two parts can be used to obtain the given Answer to the given Question. If not, extract all the necessary Context from the Question that makes it enough to obtain the given Answer to the given Question.*
- 5. ****Determine Outcome:**** Based on the reasoning, conclude and extract the Knowledge Keyword, the Attribute and the Context (if necessary) of the Question according to the Question corpus.*

Output Format

Provide the outcome in the following format:

- ****Step-by-Step Reasoning:**** [Detailed reasoning here]*
- ****Knowledge Keyword:**** [Extracted Knowledge Keyword here]*
- ****Attribute:**** [Extracted Attribute of the Knowledge Keyword here]*
- ****Context:**** [Extracted Context within the Question to make up for the Knowledge Keyword and the Attribute here if necessary]*

Examples

[examples]

Notes

- Strictly follow the format of the examples and give Knowledge Keywords, the Attribute and the Context (if necessary) anyway.*
- The extracted Knowledge Keyword, Attribute and Context (if necessary) should be the original text within the Question and should not incorporate any phrases that cannot be exactly matched in the Question.*
- Never include any information from the options of the multiple choice question, especially the content of the answer option.*
- The extracted Knowledge Keyword, Attribute and Context (if necessary) should include all the necessary information only within the Question Corpus for answering the Question and obtaining the given Answer.*

*****Fact:**** [question] [option content list] [subject] [answer option index][answer option ID]*

*****Question:**** [question]*

*****Answer:**** [content of the answer option]*

918
919
920
921
922
923
924
925
926
927
928
929
930
931
932
933
934
935
936
937
938
939
940
941
942
943
944
945
946
947
948
949
950
951
952
953
954
955
956
957
958
959
960
961
962
963
964
965
966
967
968
969
970
971

Prompt template for triplet validation & reflection:

You are an expert of [subject] and an advanced reasoning agent that can determine whether the given Knowledge Keyword, Attribute of the Knowledge Keyword and the Context present most of the necessary information of the Question for obtaining the given Answer. Suppose you have sufficient background knowledge about subj. Consider the given Knowledge Keyword, Attribute and the Context, then determine whether the given Answer can be directly obtained from them even without the Question.

Steps

1. ****Check the Semantic completeness:**** *Suppose you have sufficient background knowledge about [subject], and you can solve the given Question and obtain the given Answer. Read through the given Knowledge Keyword, Attribute, Context and the given Question. Check if the given Knowledge Keyword, Attribute, Context are the original text within the Question and contain the necessary queried information the Question itself provided (ignore the information the Question did not provided). If not so, check if the missed information is indeed incorporated in the Question (which is not acceptable, but if not, it is acceptable). Point out the information that is within the Question but they have missed. Then in a few sentences, diagnose the possible reason for failure or the phrasing discrepancy, and devise new, concise, high-level improvement suggestions to avoid the same failure.*

2. ****Check the Answer relevance:**** *Suppose you have sufficient background knowledge about subj, and you can solve the given Question and obtain the given Answer. Read through the given Knowledge Keyword, Attribute, Context and the given Question. Read through the given Knowledge Keyword, Attribute, Context and the given Answer. Check if the Answer can be directly inferred with the given Knowledge Keyword, Attribute and the Context without seeing the Question. If not so, check if the missed information is indeed incorporated in the Question (which is not acceptable, but if not, it is acceptable). Point out the information that is within the Question but they have missed. Then in a few sentences, diagnose the possible reason for failure or the phrasing discrepancy, and devise new, concise, high-level improvement suggestions to avoid the same failure.*

3. ****Check the Semantic Redundancy:**** *Read through the given Knowledge Keyword, Attribute, Context, the given Question and the given corresponding Answer. Check if the Answer can be directly matched within the given Knowledge Keyword, Attribute and the Context. Check if there are any unnecessary information within the given Knowledge Keyword, Attribute and the Context for obtaining the given Answer to the Question. If not so, point out what is redundant. Then in a few sentences, diagnose the possible reason for failure or the phrasing discrepancy, and devise new, concise, high-level improvement suggestions to avoid the same failure.*

Output Format

Provide the outcome in the following format:

- ****Step-by-Step Reasoning:**** *[Detailed reasoning here]*
 - ****Verdict for the given Knowledge Keyword, Attribute and Context:**** *[Single verdict (Yes/No) here for whether the given Knowledge Keyword, Attribute and Context contain most of the asked information of the Question, can be used to infer the given Answer with only them without the whole Question, and do not contain redundant information for obtaining the given Answer.]*

Notes

- *Do not deviate from the specified format. Do not generate anything else after the Verdict (only Yes/No) for the given Knowledge Keyword, Attribute and Context.*
 - *Suppose you have sufficient background knowledge about subj, and you can solve the given Question and obtain the given Answer. For Semantic completeness and Answer relevance, it is acceptable to miss information that is also not incorporated in the Question.*
 - *Provide a detailed explanation following the given steps before arriving at the verdict (Yes/No). Provide a final verdict (only Yes/No) in order at the end in the given format.*

- ****Question:**** *[question]*
 - ****Answer:**** *[answer]*

- ****Knowledge Keyword:**** *[extracted knowledge entity]*
 - ****Attribute:**** *[extracted attribute]*
 - ****Context:**** *[extracted context]*

972
973
974
975
976
977
978
979
980
981
982
983
984
985
986
987
988
989
990
991
992
993
994
995
996
997
998
999
1000
1001
1002
1003
1004
1005
1006
1007
1008
1009
1010
1011
1012
1013
1014
1015
1016
1017
1018
1019
1020
1021
1022
1023
1024
1025

Prompt template for the second round triplet extraction:

You are an advanced reasoning agent that can improve through self-reflection and an expert of Knowledge Keyword extraction. Analyze and summarize the Question based on the given Fact corpus and extract the Knowledge Keyword, the Attribute and the Context (if necessary) within the Question.

Given a Fact corpus, a Question about the Fact corpus, and the Answer to the Question, analyze the Question corpus as well as the given Answer. Applying the provided steps, extract the Knowledge Keyword, the Attribute of the Knowledge Keyword and the necessary Context to rephrase the Question, ensuring they are sufficient for answering the given Question and obtaining the given Answer.

Steps

1. **Review the Fact corpus:** Read through the entire Fact corpus to understand the context.
2. **Identify the Question:** Focus on the given Question to capture which part of the Fact corpus it is asking about.
3. **Understand the Answer to the Question:** Compare the given Answer and the identified questioned part within the Fact corpus and understand why this answer was chosen.
4. **Write Step-by-Step Reasoning:**
 - Identify the asked Knowledge Keyword in the Question that is the subject of the most information in the Fact corpus and the asked Question is about the information among.
 - Determine the asked Attribute of the Knowledge Keyword in the Question, which can be used to infer the given Answer.
 - Review the identified Knowledge Keyword and Attribute to confirm that only these two parts can be used to obtain the given Answer to the given Question. If not, extract all the necessary Context from the Question that makes it enough to obtain the given Answer to the given Question.
5. **Determine Outcome:** Based on the reasoning, conclude and extract the Knowledge Keyword, the Attribute and the Context (if necessary) of the Question according to the Question corpus.

Output Format

Provide the outcome in the following format:

- **Step-by-Step Reasoning:** [Detailed reasoning here]
- **Knowledge Keyword:** [Extracted Knowledge Keyword here]
- **Attribute:** [Extracted Attribute of the Knowledge Keyword here]
- **Context:** [Extracted Context within the Question to make up for the Knowledge Keyword and the Attribute here if necessary]

Examples

[examples]

You will be given a previous trial. You were unsuccessful in extracting the Knowledge Keyword, Attribute and the necessary that meet the requirements in the previous trial. Given the Reflection below, improve the process. The process is as follows:

Previous returns:

- **Fact:** [question] [option content list] [subject] [answer option index][answer option ID]
- **Question:** [question]
- **Answer:** [answer option content]
- **Knowledge Keyword:** [extracted knowledge entity of the last trial]
- **Attribute:** [attribute of the last trial]
- **Context:** [context of the last trial]
- **Reflection:** [rational of the last trial]

Notes

- Consider the Reflection given above. Improve the extraction of Knowledge Keyword, Attribute and Context (if necessary).
- Strictly follow the format of the examples and give Knowledge Keywords, the Attribute and the Context (if necessary) anyway.
- The extracted Knowledge Keyword should be phrases within the Question and should not incorporate any information of the Fact corpus or the given Answer that is not mentioned in the Question.
- The extracted Attribute and Context (if necessary) should only include information from the Question corpus. Never include information from the options of the multiple choice question, especially the content of the answer option.
- The extracted Knowledge Keyword, Attribute and Context (if necessary) should include all the necessary information only within the Question Corpus for answering the Question and obtaining the given Answer.

Fact: [question] [option content list] [subject] [answer option index][answer option ID]

Question: [question]

Answer: [content of the answer option]

1026		
1027		
1028	Original MCQ	Original MCQ Example
1029	<i>You are an expert on multiple choice questions of</i>	<i>You are an expert on multiple choice questions of</i>
1030	<i>[subject]. Analyze the given question and the given</i>	<i>high school computer science. Analyze the given</i>
1031	<i>options. Determine the correct answer option to the</i>	<i>question and the given options. Determine the cor-</i>
1032	<i>question.</i>	<i>rect answer option to the question.</i>
1033	<i>Given a Question and the potential Answer options</i>	<i>Given a Question and the potential Answer options</i>
1034	<i>to the Question, analyze the Question as well as the</i>	<i>to the Question, analyze the Question as well as the</i>
1035	<i>given options. Generate the option ID of the correct</i>	<i>given options. Generate the option ID of the correct</i>
1036	<i>option (answer).</i>	<i>option (answer).</i>
1037	- Question:	- Question:
1038	<i>[question]</i>	<i>Which of the following is usually NOT represented</i>
1039		<i>in a subroutine's activation record frame for a stack-</i>
1040	- Options:	<i>based programming language?</i>
1041	A. <i>[option A]</i>	- Options:
1042	B. <i>[option B]</i>	A. <i>Values of local variables</i>
1043	C. <i>[option C]</i>	B. <i>A heap area</i>
1044	D. <i>[option D]</i>	C. <i>The return address</i>
1045		D. <i>Stack pointer for the calling activation record</i>
1046	TrinEval MCQ	TrinEval MCQ Example
1047	<i>You are an expert on multiple choice questions of</i>	<i>You are an expert on multiple choice questions of</i>
1048	<i>[subject]. Analyze the given Knowledge Entity, At-</i>	<i>high school computer science. Analyze the given</i>
1049	<i>tribute of the Knowledge Entity, the Context of a</i>	<i>Knowledge Entity, Attribute of the Knowledge En-</i>
1050	<i>question, and the given options to the question. De-</i>	<i>ity, the Context of a question, and the given options</i>
1051	<i>termine the correct answer option to the question.</i>	<i>to the question. Determine the correct answer op-</i>
1052		<i>tion to the question.</i>
1053	<i>The Knowledge Entity is the questioned subject of</i>	<i>The Knowledge Entity is the questioned subject of</i>
1054	<i>the question. The Attribute is the questioned at-</i>	<i>the question. The Attribute is the questioned at-</i>
1055	<i>tribute of the Knowledge Entity, and the Context</i>	<i>tribute of the Knowledge Entity, and the Context</i>
1056	<i>is the necessary context information for answering</i>	<i>is the necessary context information for answering</i>
1057	<i>the question. Given a set of Knowledge Entity, At-</i>	<i>the question. Given a set of Knowledge Entity, At-</i>
1058	<i>tribute, and Context (which three are extracted as</i>	<i>tribute, and Context (which three are extracted as</i>
1059	<i>the key information from a question), and the po-</i>	<i>the key information from a question), and the po-</i>
1060	<i>tential Answer options to the Question, analyze the</i>	<i>tential Answer options to the Question, analyze the</i>
1061	<i>given Knowledge Entity, Attribute, Context as well</i>	<i>given Knowledge Entity, Attribute, Context as well</i>
1062	<i>as the options. Generate the option ID of the cor-</i>	<i>as the options. Generate the option ID of the cor-</i>
1063	<i>rect option (answer).</i>	<i>rect option (answer).</i>
1064	- Knowledge Entity:	- Knowledge Entity:
1065	<i>[knowledge entity]</i>	<i>subroutine's activation record frame</i>
1066		- Attribute:
1067	- Attribute:	<i>usually NOT represented</i>
1068	<i>[attribute]</i>	- Context:
1069		<i>for a stack-based programming language</i>
1070	- Options:	- Options:
1071	A. <i>[option A]</i>	A. <i>Values of local variables</i>
1072	B. <i>[option B]</i>	B. <i>A heap area</i>
1073	C. <i>[option C]</i>	C. <i>The return address</i>
1074	D. <i>[option D]</i>	D. <i>Stack pointer for the calling activation record</i>
1075		

Table 2: Template and an example of the Original MCQ template and the TrinEval MCQ template. [·] refers to the blank that should be filled according to the content of each MCQ.

1076

1077

1078

1079

1080
1081
1082
1083
1084
1085
1086
1087
1088
1089
1090
1091
1092
1093
1094
1095
1096
1097
1098
1099
1100
1101
1102
1103
1104
1105
1106
1107
1108
1109
1110
1111
1112
1113
1114
1115
1116
1117
1118
1119
1120
1121
1122
1123
1124
1125
1126
1127
1128
1129
1130
1131
1132
1133

LLMs	Dataset	2 × 2 squares		3 × 3 squares	
		Ratio (%)	Pearson correlation	Ratio (%)	Pearson correlation
Llama2-Qwen	All	38.57	-0.7755	74.17	-0.8124
	Simple	37.07	-0.7784	72.63	-0.8121
	Pro	38.66	-0.783	74.51	-0.8109
Llama2-GPT	All	35.22	-0.7835	71.04	-0.7924
	Simple	33.9	-0.777	69.62	-0.7919
	Pro	35.45	-0.7916	71.54	-0.7881
Mistral-Qwen	All	44.9	-0.8722	80.82	-0.8794
	Simple	38.47	-0.8494	74.04	-0.8271
	Pro	44.32	-0.8045	80.08	-0.8682
Mistral-GPT	All	40.37	-0.8042	76.58	-0.8736
	Simple	35.51	-0.8297	72.27	-0.8664
	Pro	38.52	-0.7103	74.91	-0.7969
Vicuna-Qwen	All	42.94	-0.8771	79.23	-0.8365
	Simple	37.86	-0.758	73.85	-0.7168
	Pro	42.01	-0.8609	77.86	-0.886
Vicuna-GPT	All	38.69	-0.8621	74.83	-0.8672
	Simple	34.77	-0.8096	70.71	-0.7775
	Pro	37.37	-0.7794	73.98	-0.8728

Table 3: The ratio and the Pearson-correlation between the F_c and F_m of the MCQs within the upper right and lower left 2×2 and 3×3 squares. For LLMs, ‘Llama2-Qwen’ refers that the F_c and F_m are calculated with Llama2 based on the Qwen-extracted triplet, and similarly hereinafter. For the Dataset column, ‘All’ stands for all the qualified MCQs after the triplet extraction, ‘Pro’ refers to the qualified MCQs that are the members of the mmlupro dataset while ‘Simple’ refers to the rest of the MCQs that are relatively easier.

Model	TrinEval correct only	both correct	original correct only	K.P. score
Qwen	0.0%	96.296%	3.704%	4.101
GPT	0.0%	96.667%	3.333%	4.369

Table 4: Result of Human Annotation on LLM-based knowledge preserving (K.P.) evaluation in TrinEval

Threshold	1%		3%		5%	
Subset	RM	GCL	RM	GCL	RM	GCL
Llama-Qwen	29.538	16.501	32.718	17.281	33.790	18.060
Llama-GPT	28.925	17.777	31.327	18.275	32.315	19.068
Mistral-Qwen	78.663	47.648	86.003	52.995	89.598	55.525
Mistral-GPT	82.932	56.375	90.597	62.112	94.384	64.820
Vicuna-Qwen	38.761	21.006	41.396	23.189	42.490	24.185
Vicuna-GPT	37.037	22.996	41.396	23.189	42.490	24.185

Table 5: Result of averaged embedding distance of the closest MCQs under different thresholds