# RECKON: Large-scale <u>Rec</u>erence-based Efficient <u>K</u>nowledge Evaluati<u>on</u> for Large Language Model

**Lin Zhang**♠*, **Zhouhong Gu**♠*, **Xiaoran Shi**♠, **Hongwei Feng**♠† , **Yanghua Xiao**♠†

♠Shanghai Key Laboratory of Data Science, School of Computer Science, Fudan University

{linzhang22, zhgu22, xrshi21}@m.fudan.edu.cn

{hwfeng, shawyh}@fudan.edu.cn

## Abstract

As large language models (LLMs) advance, efficient knowledge evaluation becomes crucial to verifying their capabilities. Traditional methods, relying on benchmarks, face limitations such as high resource costs and information loss. We propose the Large-scale Reference-based Efficient Knowledge Evaluation for Large Language Model (RECKON), which directly uses reference data to evaluate models. RECKON organizes unstructured data into manageable units and generates targeted questions for each cluster, improving evaluation accuracy and efficiency. Experimental results show that RECKON reduces resource consumption by 56.5% compared to traditional methods while achieving over 97% accuracy across various domains, including world knowledge, code, legal, and biomedical datasets. Code is available at https://github.com/MikeGu721/reckon.

## 1 Introduction

As large language models (LLMs) continue to advance rapidly, knowledge evaluation has become an essential component for verifying their capabilities and driving continuous improvements. The fundamental objective of evaluation is to assess and ensure the LLM's alignment with reference data. Thorough and accurate knowledge evaluation plays a critical role in ensuring reliable model performance across different fields, especially in crucial areas such as healthcare, education, and law, where strong evaluation methods are necessary to reduce risks from errors or outdated information. With the rapid growth of LLM applications, the requirements for knowledge evaluation have become increasingly extensive. Modern evaluation approaches must meet growing demands, including wider coverage of different fields, more complex analytical tasks, and the handling of much larger
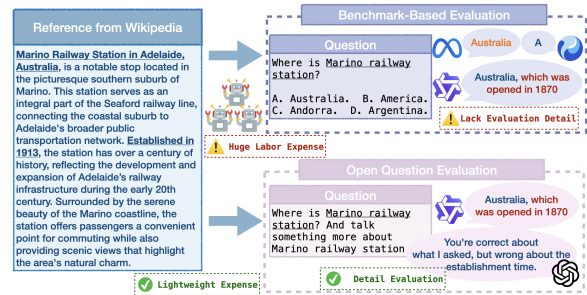


Figure 1: The comparison of different evaluation method.

datasets. These requirements pose significant challenges to traditional evaluation methods.

Traditional knowledge evaluation methodologies primarily rely on benchmarks, encompassing various assessment formats such as multiple-choice questions (Gu et al., 2024), fill-in-the-blanks exercises (Hendrycks et al., 2020), and marking tasks (Dimitrov et al., 2024). These benchmarks fundamentally serve as condensed or refined representations of reference data, carefully designed to facilitate manageable and systematic evaluation processes. Despite widespread adoption across the field, these benchmark-based approaches exhibit several inherent limitations. First, the creation and verification of comprehensive benchmarks demand substantial resources, necessitating significant manual and computational expenses. Second, the process of distilling complex, multifaceted information into standardized questions inevitably results in information loss, potentially introducing systematic biases and compromising the validity of knowledge assessment. Finally, individual benchmarks face considerable challenges in fully capturing the intricate complexity and expansive diversity of knowledge domains, thereby constraining their adaptability and broader applicability.

As illustrated in Fig.1, directly leveraging rich and dynamic reference data is an alternative ap-

proach, which presents distinct advantages over traditional benchmark-based methods. This direct reference approach not only preserves the nuanced complexity of the original information but also enables more comprehensive and authentic evaluation scenarios. Furthermore, it facilitates dynamic updating of evaluation criteria as new knowledge emerges, ensuring continued relevance and accuracy. The approach also demonstrates superior scalability across different knowledge domains and evaluation contexts, offering enhanced flexibility in evaluation. However, this direct reference-based evaluation introduces its own set of significant challenges: First, reference datasets typically exist in large-scale, unstructured formats, presenting substantial difficulties in effective organization and utilization during evaluation processes. Second, the absence of predefined questions or standardized answer formats complicates the establishment of consistent and equitable alignment between model outputs and reference data.

To address these challenges, we propose a <u>R</u>eference-based <u>E</u>fficient <u>K</u>nowledge Evaluati<u>on</u> framework, called RECKON, which directly and efficiently evaluates LLMs using reference data. In response to the difficulty of organizing and utilizing large and unstructured reference datasets, RECKON decomposes the reference data into discrete, manageable knowledge units. These units are then grouped into thematic clusters, forming an organized knowledge structure that streamlines the evaluation process. To solve the problem of ensuring consistent and fair alignment without predefined questions or fixed answers, RECKON generates targeted questions for each knowledge cluster. The evaluation focuses on assessing whether the model's responses adequately cover all relevant knowledge units, ensuring comprehensiveness and accuracy.

Experimental results highlight the superiority of RECKON as a comprehensive, efficient, and adaptable framework for knowledge evaluation. Across various common-sense and anti-common-sense tasks, RECKON effectively adapts references to reduce the inherent biases of LLMs. Notably, RECKON reduced resource consumption by 56.5% compared to "full reference input evaluation" without compromising evaluation accuracy. Additionally, RECKON demonstrated exceptional adaptability across diverse domains, including world knowledge, code, legal, and biomedical datasets, achieving an accuracy exceeding 97%.

## 2 Related work

Open-ended question answering benchmarks for models include MMLU (Hendrycks et al., 2020), C-eval (Huang et al., 2024), Xiezhi (Gu et al., 2024), etc. These benchmarks contain a series of questions that are described in natural language and require the model to give an open-ended answer. However, these benchmarks rely on a large number of manual annotations and cannot be updated with the latest knowledge. Currently, some methods propose using LLMs to automatically build updatable benchmarks, such as LM-as-an-Examiner (Bai et al., 2024) and TreeEval (Li et al., 2024). In these methods, when no benchmark is available, LLMs generate questions based on their own knowledge. However, the knowledge in LLM cannot be inherently complete, and the language model is biased, which will lead to incomplete questions and deviation of questions to a certain extent.

Traditional open evaluation metrics are based on n-grams to measure semantic similarity between texts, including BLEU (Papineni et al., 2002), ROUGE (Lin, 2004), METEOR (Banerjee and Lavie, 2005), etc. These methods enable the automated evaluation of natural language texts. However, these methods lack expressive power and cannot semantically distinguish key information, such as negative words. To address these problems, BERTScore (Zhang et al., 2019), GPTScore (Fu et al., 2023), and other methods use language models to derive similarities between candidate answers and reference answers. In the latest evaluation work on LLM, LLM-as-Judge (a paradigm that uses LLM to evaluate open results) (Zheng et al., 2024; Chan et al., 2023) is proposed. This method is recognized for its strong interpretability and scalability, making it widely adopted. These methods still rely on either artificially constructed reference answers or the LLM's own knowledge, leaving the issue of knowledge limitations unresolved.

RECKON distinguishes itself from traditional benchmarks by using a dynamic, reference-based approach rather than static question-answer pairs. This design makes it more adaptable, scalable, and effective at evaluating nuanced knowledge across diverse domains. Unlike traditional benchmarks, which are limited in coverage, flexibility, and efficiency, RECKON leverages rich reference data to align with evolving knowledge. Moreover, RECKON distinguishes itself from other LLM-related evaluation methods, such as LLM-as-Judge,

by grounding evaluations in external references rather than relying on the model's subjective judgment. This ensures a more objective, transparent, and consistent evaluation process, effectively mitigating biases and variability inherent in LLM-as-Judge approaches.

## 3  RECKON Framework

As shown in Fig.2, RECKON consists of five main components, and the description of these components is as follows:

**Step 1 Preparing Knowledge Unit Candidates:** The RECKON framework begins by preparing knowledge units that will be utilized throughout the evaluation process. These units are sourced primarily from two places: External references, such as curated datasets, textbooks, or knowledge bases relevant to the target domain. And Previously evaluated knowledge units, units from earlier rounds of assessment that were flagged as uncovered or incomplete by the evaluated model.

**Step 2 Clustering:** After compiling the knowledge units, clustering step involves organizing them into meaningful clusters to streamline the evaluation process. Clustering effectively organizes knowledge units into coherent groups to facilitate the evaluation process. This organization is achieved by analyzing the semantic relationships between different knowledge units. Each knowledge unit, which may consist of facts, concepts, and information relevant to the related field, is transformed into an embedding representation. These units are then classified based on the proximity of their embeddings in semantic space, effectively creating subsets of knowledge that are thematically related.

**Step 3 Question Generation:** In this step, RECKON employs a LLM to generate targeted questions for each cluster of knowledge units. This step utilizes the capabilities of the language model to pose questions. The questions are designed to align with the central theme of the cluster, ensuring they are contextually appropriate and address the core knowledge units. These questions are tailored to test not only surface-level understanding but also deeper insights, relationships, and nuances within the knowledge units. During the evaluation, if the language model's responses reveal gaps, ambiguities, or misconceptions, RECKON dynamically adjusts the questions to probe these areas further. This iterative refinement ensures a comprehensive assessment of the evaluated model's knowledge and reasoning capabilities.

**Step 4 Get Response:** Once the questions are generated, they are presented to get responses from the evaluated targets. The responses obtained during this phase serve as the primary input for analysis. RECKON evaluates the accuracy, completeness, and relevance of the responses in relation to the knowledge units. If a response fails to address a question, RECKON adequately will continue to prompt the model to generate questions to cover the relevant knowledge units. This ensures that every relevant knowledge unit within the cluster is thoroughly examined.

**Step 5 Judging:** This step involves instructing the evaluation LLM to assess the extent to which the evaluated model's responses cover the specified knowledge units. In this phase, RECKON labels the covered knowledge units as either correct or incorrect. For knowledge units that remain uncovered by the evaluated model's response, RECKON identifies and retains them for reorganization and re-evaluation. The process starts again from the initial step in subsequent assessment rounds.

Further technical details, including the prompts, models, inputs, and outputs used at each step, can be found in Appendix A.1.

## 4  Experiment Setup

### 4.1  Dataset and models

Our experiments encompassed diverse textual corpora spanning multiple domains:

- **Wikitext** from **Wiki** (Vrandečić and Krötzsch, 2014), comprising meticulously curated articles encompassing comprehensive knowledge. The diversity and comprehensiveness of Wikitext make it an excellent benchmark for assessing the breadth and depth of a model's general knowledge capabilities.

- **Code** from **CodeGPT** (Xiaoxuan et al., 2023), containing programming-oriented assignments. It is specifically used to evaluate the models' ability to understand and generate code or solve programming challenges.

- **Legal.term** from **Legalbench** (Guha et al., 2024), incorporating juridical definitions emphasizing legislative interpretation. It is a resource for testing a model's performance in legal reasoning, terminology comprehension, and domain-specific text generation.
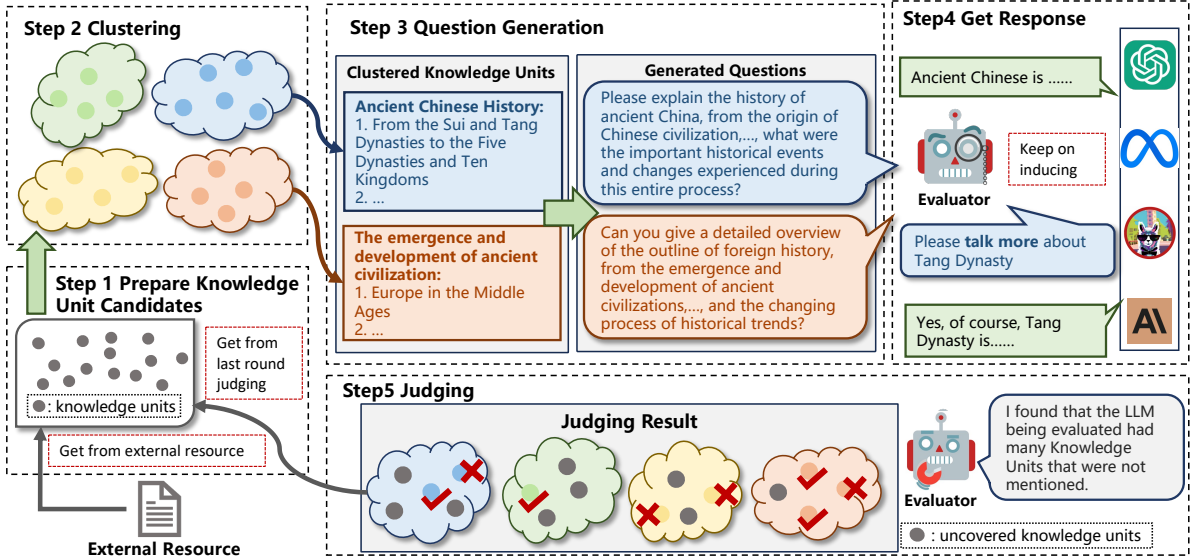
Figure 2: RECKON consists of five main components that form a cycle of assessment.

- **Med.rand** and **Med.sim** from **Pub-MedQA** (Jin et al., 2019). Med.rand includes randomized selections from the PubMedQA dataset to evaluate general biomedical knowledge. Med.sim is a similarity-based subset specifically designed to test the models' ability to identify nuanced relationships between biomedical entities.

- **Uncommon** from **Feverous** (Aly et al., 2021), encompassing descriptions that deliberately deviate from established factual paradigms. It assesses a model's ability to handle the contradiction between an evaluation model and the external reference.

Details of these datasets is presented in Tab.1.

The models under evaluation include the widely used **GPT4-turbo** (GPT4), **GPT3.5-turbo** (GPT3.5), and **LLaMA2-chat-13b** (LLaMA2). Additionally, we incorporated the **PMC-LLaMA-13b**(PMC), fine-tuned on the LLaMA2 architecture, specifically tailored for the biomedical field, to assess its performance within the related domain.

### 4.2 Baseline

In our evaluation, we compare RECKON against several baseline methods: **BLEU**, a probabilistic measure based on n-gram matching by comparing the overlap of n-grams between two sentences. **GPT Score**, an embedding-based indicator utilizing pre-trained language model embeddings to compare the semantic similarity between generated and reference texts. **LLM-as-Judge**

| Dataset | Source | #Words | #Para. |
|---------|--------|--------|--------|
| Code | CodeGPT[1] | 2,371k | 5,511 |
| Legal.term | Legalbench[2] | 233k | 695 |
| Wikitext | Wiki[3] | 11,834k | 4,396 |
| Uncommon | Feverous[4] | 49k | 2,000 |
| Med.rand | PubMedQA[5] | 877k | 2,000 |
| Med.sim | PubMedQA | 351k | 2,000 |

Table 1: Detail statistic information of the datasets encompassed in the experiments.

**w/o reference(Jw/oR)**, adopting the LLM-as-judge paradigm and relying on the evaluator (LLM) to judge response correctness. **LLM-as-Judge w/ reference(Jw/R)**, assessing response correctness by aligning the answer provided by the model and the related textual context in the datasets one-by-one, which is superior to RECKON with more detail.

To establish a ground truth for comparison, we conducted **Human** evaluations as detailed in Appendix A.3. Three expert annotators assessed model responses across all datasets with access to the original reference texts. Each response was labeled as either correct or incorrect, based on its alignment with the reference text and overall quality. The human evaluations served as the gold standard for our study, allowing us to measure the accuracy and reliability of RECKON and other baseline methods. The average human scores and inter-

---

[1] https://github.com/zxx000728/CodeGPT
[2] https://hazyresearch.stanford.edu/legalbench/
[3] https://huggingface.co/datasets/wikitext
[4] https://github.com/Raldir/FEVEROUS
[5] https://pubmedqa.github.io

annotator agreement are presented in A.3.

## 5 Performance

Tab.2 displays the performance of all baselines measured by their correlation with human annotations. Tab.4 displays the resource consumption of two strong baselines, Jw/R and RECKON. The conclusions from the results are as follows:

**LLM-based methods exhibit significantly higher evaluation accuracy compared to both n-gram-based and embedding-based approaches.** The correlation with human baselines has notably improved, ranging from approximately 0.3 for BLEU to 0.6 for GPT Score and approaching nearly 1 for Jw/R, Jw/oR, and RECKON. This substantial enhancement underscores the superior efficacy of LLM-based techniques in assessing model responses in alignment with human expectations. The inferior performance of n-gram-based and embedding-based methods can be attributed to their limited understanding of the text's intrinsic meaning. In contrast, LLMs offer robust capabilities in comprehending textual context, thereby enhancing the effectiveness of Jw/oR methodologies.

**External references enhance evaluation accuracy and stability by countering internal model biases.** Firstly, references improve the accuracy of evaluation. Both the Jw/R and RECKON baselines demonstrate a higher correlation, utilizing textual references to augment the LLM's precision in assessing answer correctness. External references effectively counteract the potential influence of erroneous internal model knowledge, as previously investigated in the literature (Xie et al., 2023). Secondly, adopting external references leads to superior stability and robustness against LLM alone. Notably, the Jw/R and RECKON methods exhibit significantly lower variances of 1.9e-5 and 2.5e-4, respectively. In contrast, the Jw/oR baseline demonstrates a higher variance of 0.04, indicating greater susceptibility to internal model biases.

**RECKON reduces the cost of evaluation and maintains competitive correlation with Jw/R.** Tab.4 demonstrates the efficiency advantage of RECKON. In contrast to Jw/R, which necessitates more frequent API calls and processes greater data volumes, thereby escalating operational expenses, RECKON optimizes these aspects by concentrating on crucial tokens during the Judgment

| | GPT4 | GPT3.5 | LLaMA2 |
|---|---|---|---|
| *N-Gram Based* | | | |
| **BLEU** | 28.1 | 38.4 | 9.4 |
| *Embedding Based* | | | |
| **GPTScore** | 60.6 | 62.3 | 41.0 |
| *LLM-as-Judge Based* | | | |
| **Jw/R** | **99.8** | 98.9 | **99.2** |
| **Jw/oR** | 88.9 | 98.4 | 86.1 |
| **RECKON** | 97.6 | **99.5** | 96.3 |

Table 2: Correlations between different methods and human baseline, where the bold font indicates the highest correlation.

---

**Question:** What was the Silent Holy Stones' director award in 2005?
**Answer:** Directed by Pema Tseden, won the Golden Rooster Award for Best Directorial Debut in 2005, marking a significant step for Tibetan cinema.
**Judgement:** Incorrect

**Reference Text:** The Silent Holy Stones was nominated in the **Best Director category in the 2005 Beijing College Student Film Festival**.

---

Table 3: Performance of GPT-4 on Uncommon Dataset by RECKON.

and Question generation phases. This strategy reduces unnecessary overhead and lowers API call rates, enabling cost savings and effective evaluation. A deeper analysis in A.4 shows that method's cost is related to model capacity, with higher-capacity models incurring lower expenses compared to Jw/R.

## 6 Analysis

### 6.1 Evaluation Results

As shown in Tab.5, **LLMs demonstrate superior performance across diverse knowledge domains, with GPT-4 consistently achieving the highest accuracy scores among all tested models.** GPT-4 achieves remarkable accuracy rates of 95.4% on Code, 94.6% on Legal.term, and 97.9% on Med.sim datasets. These outstanding results underscore its versatility in handling both general and specialized tasks, particularly in biomedical domains. When compared to GPT-3.5, which achieves 90.8% on Legal.term and 54.5% on Wiki, the performance gap clearly illustrates the advantages of larger-scale models in knowledge-intensive tasks.

As shown in Tab.5, **domain-specific fine-tuning yields substantial improvements in model performance within targeted fields.** This is evidenced by PMC-LLaMA's performance in the biomedical

| | Eval. Target | #Query-API | #Token | #Judge-API | #Token | #Money |
|---|---|---|---|---|---|---|
| **Jw/R** | GPT3.5 | 1,423 | 650,098 | 1,423 | 482,779 | 38.14 |
| **RECKON** | GPT3.5 | 256 | 26,880 | 1,622 | 616,760 | 2.44 |
| **RECKON** | GPT4 | 244 | 21,692 | 1,593 | 556,257 | 2.39 |
| **RECKON** | LLaMA | 269 | 69,271 | 1,850 | 660,350 | 3.27 |

Table 4: Metrics of resource consumption measured during the Wiki evaluation phase. **Eval. Target** signifies the language model subjected to assessment. **#Query-API** and **#Judge-API** enumerate the aggregate API invocations throughout the question formulation and assessment phases, respectively. **#Token** encompasses the collective token consumption across all operations. **#Money** encapsulates the overall monetary expenditure incurred through API utilization.
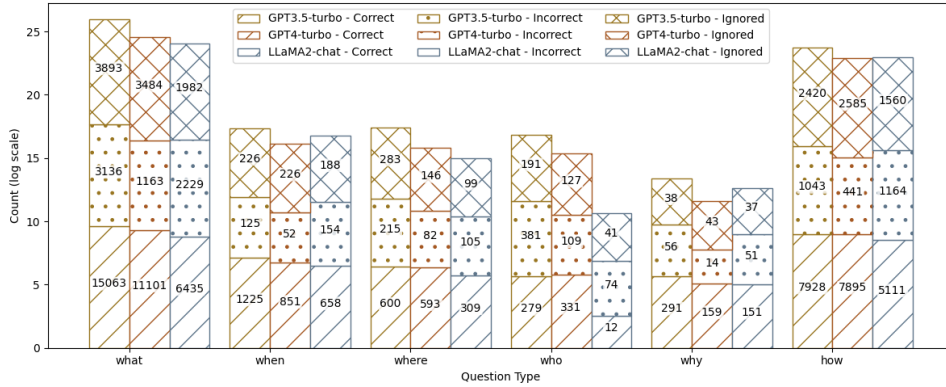


Figure 3: The counts of 6 types of questions of RECKON used for different models.

| Model | GPT4 | GPT3.5 | LLaMA2 | PMC |
|---|---|---|---|---|
| **Code** | **95.4** | 92.9 | 83.1 | – |
| **Legal.term** | **94.6** | 90.8 | 81.2 | – |
| **Wiki** | **82.6** | 54.5 | 42.9 | – |
| **Med.rand** | **95.3** | 82.1 | 73.2 | 78.5 |
| **Med.sim** | **97.9** | 82.3 | 72.5 | 75.7 |

Table 5: The performance of the model on each dataset, expressed as accuracy, where the bold font indicates the highest accuracy.

| | Code | Legal.term | Wiki | Med.rand |
|---|---|---|---|---|
| **Jw/oR** | 97.2 | 95.5 | 72.6 | 90.9 |
| **RECKON** | 95.4 | 90.9 | 54.5 | 82.1 |

Table 6: Evaluation result of RECKON and raw LLM-as-a-Judge.

| Model Accuracy | | | |
|---|---|---|---|
| | GPT4 | GPT3.5 | LLaMA2 |
| **RECKON** | 31.7 | 32.8 | 24.7 |
| **Jw/oR** | 84.4 | 87.3 | 67.5 |

| Correlation Between Human Annotation | | | |
|---|---|---|---|
| | GPT4 | GPT3.5 | LLaMA2 |
| **RECKON** | 90.57 | 92.45 | 96.22 |

Table 7: Evaluation performance on UnCommon.

## 6.2 Evaluation Performance with Reference Data

As shown in Tab.6, **models exhibit a bias toward providing 'correct' judgments when evaluating without external references.** While Jw/oR assigns consistently higher scores to GPT-3.5 compared to RECKON, manual verification reveals that RECKON actually achieves higher evaluation accuracy. This discrepancy suggests that Jw/oR's inflated results stem from its tendency to classify wrong answers as correct incorrectly.

As shown in Tab.3, **reference-based evaluation through RECKON ensures objective assessment by strictly adhering to provided reference ma-**

domain, where it achieves accuracy scores of 78.5% and 75.7% on Med.rand and Med.sim, respectively. While these scores don't surpass GPT-4's performance, they represent a significant improvement over general-purpose models like LLaMA2, which only achieves 73.2% and 72.5% on the same datasets.

| Comparison of RECKON and LLM-as-Judge | | |
|---|---|---|
| **Jw/oR** | **Question:** When was Marino railway station in Adelaide opened?<br>**Answer:** Marino Railway Station was opened on October 19, 1885.<br>**Judgement:** Correct✗<br>**Reason:** The candidate's answer correctly identifies the location of Marino Railway Station and its opening date. | **Text:** Marino railway station is on the Seaford railway line , in the southern Adelaide suburb of Marino , South Australia Once the terminus of the line that now runs south to Seaford . **Marino railway station opened in 1915 when the line was built from Adelaide to Marino.** |
| **RECKON** | **Question:** What is the common trait of Marino railway station in Adelaide, . . . , and what years were they opened?<br>**Answer:** Marino Railway Station is located in Adelaide, Australia. It was opened on October 19, 1885, and it is still operational. It serves the Seaford and Tonsley lines of Adelaide's metropolitan train network.<br>**Judgement:** Incorrect✓<br>**Reason:** The candidate incorrectly states that Marino Railway Station was opened on October 19, 1885, while the reference text indicates it opened in 1915. | |
| Performance of GPT3.5-Turbo | | |
| **Jw/R** | **Question:** In legal terms, what does the term 'damage' specifically refer to?<br>**Answer:** In legal terms, the term 'damage' refers to the harm or injury caused to a person, property, or reputation as a result of the wrongful act of another. . .<br>**Judgement:** Incorrect | **Text:** The word 'damage' traditionally describes a harm to property (hence, 'property damage'), rather than harm to the person (usually referred to as 'personal injury'). |
| **RECKON** | **Question:** In area of law, . . . , and what is the distinction between 'damage' and 'injury'?<br>**Answer:** The distinction between 'damage' and 'injury' is that damage refers to harm caused to property, while injury refers to the victim by certain injury.<br>**Judgement:** Correct | |
| **Human** | **Question:** In legal terminology, what kind of object does 'damage' apply to ?<br>**Answer:** 'Damage' usually applies to damage to objects. This includes loss or damage to property, land, equipment, etc. . .<br>**Judgement:** Correct | |

Table 8: Comparison of RECKON and Jw/oR at the top of the table, and performance of GPT3.5-Turbo of different baselines at the bottom of the table.

**terials, even when they conflict with real-world knowledge.** This is exemplified in the case where GPT-4's response about a film "winning the Golden Rooster Award for Best Directorial Debut in 2005" was marked incorrect due to contradicting reference information, despite being factually accurate. In contrast, using Jw/oR, GPT-4 achieved high accuracy on the Uncommon dataset by relying on its internal knowledge. Manual verification confirms RECKON's strict adherence to reference materials in over 90% of cases, regardless of potential conflicts with the model's internal knowledge base.

## 6.3 Evaluation Performance Knowledge Clustering and Concept Characteristics

As shown in Figure 4a, **knowledge clusters with high intra-cluster similarity and low inter-cluster overlap optimize RECKON's evaluation accuracy.** The analysis reveals 13 distinct clusters characterized by high intra-cluster semantic similarity and large inter-cluster distances in the feature space. Between rounds 30 (Figure 4a) and 40 (Figure 4b), domain-specific clusters such as "Various films", "TV shows", "operas", and "related media productions and Various STEM topics" achieve complete evaluation coverage. In contrast, heterogeneous clusters containing mixed domains, like *"Various historical, biographical, electoral, media, educational and legal topics"*, exhibit high feature dispersion in the embedding space, leading to reduced evaluation precision.

As shown in Fig.5, **the hierarchical clustering strategy in RECKON optimizes evaluation efficiency through systematic knowledge unit distribution.** The frequency analysis demonstrates an inverse relationship between knowledge unit sampling frequency and the occurrence of 'ignored' units. This pattern validates that RECKON's clustering algorithm effectively prevents redundant sampling of high-frequency knowledge units while maintaining comprehensive coverage across the knowledge space.

As shown in Fig.5, **LLMs' performance exhibits systematic degradation when processing knowledge units with high information density.** This is evidenced by the increased frequency of 'ignored' knowledge units containing extensive information, such as *"List of medical drama television programs worldwide"*. The performance degradation is particularly notable in sparse clusters, suggesting that LLMs' processing efficiency decreases as the information density of the input

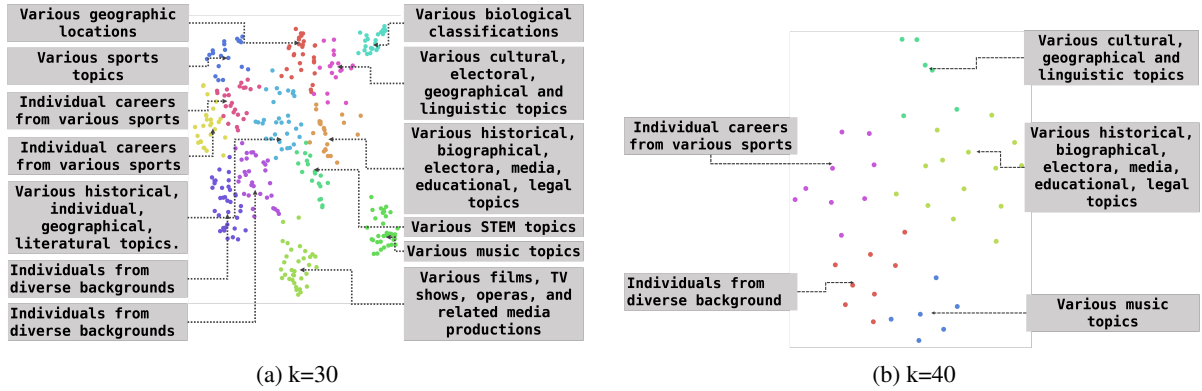7

(a) k=30             (b) k=40

Figure 4: The embedding map of the remaining knowledge units to be evaluated on the Wiki dataset in the process of evaluating GPT3.5-turbo at round **k** (k is different of the two sub-figures), where the representation of knowledge units of the same color is clustered in the same class, and their class names are displayed.
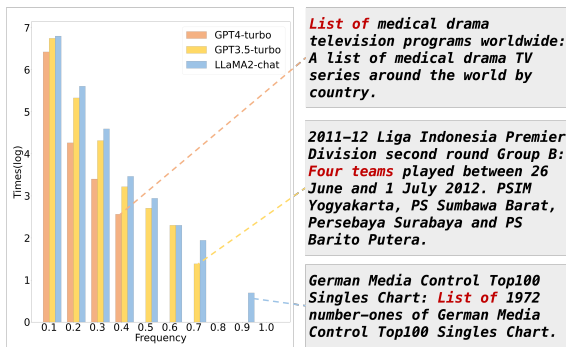


Figure 5: The frequency of knowledge units in multiple iterations of RECKON. The X-axis is the frequency of knowledge units, and the vertical axis is the number of knowledge units of a certain frequency.

content increases.

### 6.4 Evaluation Performance of Different Types of Questions

As shown in Fig.3, **LLMs consistently exhibit lower performance on identity-based and location-based queries.** The data reveals a systematic underperformance across all models when handling 'who' questions, likely due to the complexities of processing person-specific information and historical contexts. Similar performance deficits are observed with 'where' questions, suggesting challenges in managing geographical details and location-based information. This pattern indicates fundamental limitations in models' ability to process and accurately represent identity and location-specific knowledge.

As shown in Tab.8, **RECKON enhances model performance by leveraging contextual information to activate inherent reasoning capabilities.** While GPT-4-turbo and LLaMa2-chat show supe-

rior performance with Jw/R on factual recall questions ('who' and 'where'), RECKON achieves comparable or better results with models like GPT-3.5-turbo by incorporating rich contextual information. For instance, RECKON improves GPT-3.5-turbo's understanding of 'damage' concepts by providing relevant contextual cues, whereas Jw/R struggles without such supporting information.

### 6.5 Case Study

As shown in Tab.2, **RECKON demonstrates superior evaluation accuracy compared to Jw/oR through its reference-based assessment approach.** This is exemplified by the case of the Marino Railway Station opening time query, where Jw/oR incorrectly marks an answer as 'Correct' based on the model's internal knowledge, while RECKON accurately identifies it as 'Incorrect' by comparing it against the Wiki reference. This case highlights RECKON's ability to provide reliable evaluations through reference-based verification rather than depending solely on model knowledge.

### 7 Conclusion

In this paper, we introduce RECKON, a novel method for evaluating large language models across diverse domains. Leveraging LLMs' ability to understand and respond to complex queries, RECKON addresses the limitations of traditional evaluation methods. Our experiments on various datasets show that RECKON achieves high consistency with human evaluation, which highlights its effectiveness in assessing model responses. This approach not only offers a scalable and efficient means to evaluate LLMs but also advances the field of model evaluation in knowledge-intensive tasks.

## Limitation

While the proposed RECKON offers significant advancements in evaluating text generated by language models, it is not without limitations. RECKON's effectiveness heavily depends on the quality and relevance of the external reference materials used. If these references are incomplete or outdated, the framework's evaluations may be compromised. The challenge of integrating large volumes of reference data remains, as even with synthesized knowledge units, the risk of overlooking critical details or context persists. The sensitivity of LLMs to prompt variations can also result in inconsistent evaluation outcomes when dealing with diverse or ambiguously phrased questions. This variability in model responses may affect the reliability of the evaluation results, particularly in scenarios where nuanced understanding is crucial.

## Ethical Concerns

The RECKON introduces several ethical concerns. The external reference data used may include sensitive or controversial content, which could lead to the perpetuation of biases or misinformation. Additionally, handling proprietary or personal information raises privacy and intellectual property concerns.

To address these issues, we implement strict protocols to vet reference data for sensitivity and relevance. We ensure transparency in our data curation process and prioritize ethical standards to safeguard privacy and prevent misuse, balancing the benefits of comprehensive evaluation with responsible data handling.

## References

Rami Aly, Zhijiang Guo, Michael Schlichtkrull, James Thorne, Andreas Vlachos, Christos Christodoulopoulos, Oana Cocarascu, and Arpit Mittal. 2021. Feverous: Fact extraction and verification over unstructured and structured information. *arXiv preprint arXiv:2106.05707*.

Yushi Bai, Jiahao Ying, Yixin Cao, Xin Lv, Yuze He, Xiaozhi Wang, Jifan Yu, Kaisheng Zeng, Yijia Xiao, Haozhe Lyu, et al. 2024. Benchmarking foundation models with language-model-as-an-examiner. *Advances in Neural Information Processing Systems*, 36.

Satanjeev Banerjee and Alon Lavie. 2005. Meteor: An automatic metric for mt evaluation with improved correlation with human judgments. In *Proceedings of the acl workshop on intrinsic and extrinsic evaluation measures for machine translation and/or summarization*, pages 65–72.

Chi-Min Chan, Weize Chen, Yusheng Su, Jianxuan Yu, Wei Xue, Shanghang Zhang, Jie Fu, and Zhiyuan Liu. 2023. Chateval: Towards better llm-based evaluators through multi-agent debate. *arXiv preprint arXiv:2308.07201*.

Dimitar Dimitrov, Firoj Alam, Maram Hasanain, Abul Hasnat, Fabrizio Silvestri, Preslav Nakov, and Giovanni Da San Martino. 2024. Semeval-2024 task 4: Multilingual detection of persuasion techniques in memes. In *Proceedings of the 18th International Workshop on Semantic Evaluation (SemEval-2024)*, pages 2009–2026.

Jinlan Fu, See-Kiong Ng, Zhengbao Jiang, and Pengfei Liu. 2023. Gptscore: Evaluate as you desire. *arXiv preprint arXiv:2302.04166*.

Zhouhong Gu, Xiaoxuan Zhu, Haoning Ye, Lin Zhang, Jianchen Wang, Yixin Zhu, Sihang Jiang, Zhuozhi Xiong, Zihan Li, Weijie Wu, et al. 2024. Xiezhi: An ever-updating benchmark for holistic domain knowledge evaluation. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pages 18099–18107.

Neel Guha, Julian Nyarko, Daniel Ho, Christopher Ré, Adam Chilton, Alex Chohlas-Wood, Austin Peters, Brandon Waldon, Daniel Rockmore, Diego Zambrano, et al. 2024. Legalbench: A collaboratively built benchmark for measuring legal reasoning in large language models. *Advances in Neural Information Processing Systems*, 36.

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*.

Yuzhen Huang, Yuzhuo Bai, Zhihao Zhu, Junlei Zhang, Jinghan Zhang, Tangjun Su, Junteng Liu, Chuancheng Lv, Yikai Zhang, Yao Fu, et al. 2024. C-eval: A multi-level multi-discipline chinese evaluation suite for foundation models. *Advances in Neural Information Processing Systems*, 36.

Qiao Jin, Bhuwan Dhingra, Zhengping Liu, William W Cohen, and Xinghua Lu. 2019. Pubmedqa: A dataset for biomedical research question answering. *arXiv preprint arXiv:1909.06146*.

Xiang Li, Yunshi Lan, and Chao Yang. 2024. Treeeval: Benchmark-free evaluation of large language models through tree planning. *arXiv preprint arXiv:2402.13125*.

Chin-Yew Lin. 2004. Rouge: A package for automatic evaluation of summaries. In *Text summarization branches out*, pages 74–81.

Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the 40th annual meeting of the Association for Computational Linguistics*, pages 311–318.

Denny Vrandečić and Markus Krötzsch. 2014. Wikidata: a free collaborative knowledgebase. *Communications of the ACM*, 57(10):78–85.

Zhu Xiaoxuan, Xiong Zhuozhi, Zhang Lin, Ye Haoning, Gu Zhouhong, Li Zihan, Jiang Sihang, Feng Hongwei, Xiao Yanghua, Wang Zili, Yang Dongjie, and Wang Shusen. 2023. Codegpt: A code-related dialogue dataset generated by gpt and for gpt. `https://github.com/zxx000728/CodeGPT`.

Jian Xie, Kai Zhang, Jiangjie Chen, Renze Lou, and Yu Su. 2023. Adaptive chameleon or stubborn sloth: Revealing the behavior of large language models in knowledge conflicts. *arXiv preprint arXiv:2305.13300*.

Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q Weinberger, and Yoav Artzi. 2019. Bertscore: Evaluating text generation with bert. *arXiv preprint arXiv:1904.09675*.

Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric Xing, et al. 2024. Judging llm-as-a-judge with mt-bench and chatbot arena. *Advances in Neural Information Processing Systems*, 36.

## A  Appendix

### A.1  Detials of RECKON

#### A.1.1  Prompt

In Fig.6, Fig.7, and Fig.8, we illustrate three distinct prompts, each designed to complete different tasks within RECKON. Fig.6 represents a question generation prompt in step 3 of **Question Generation**. Fig.7 represents a question-answering prompt in step 4 of **Get Response**. Fig.8 represents a response judgment prompt in step 5 of **Judging**.

Each of them includes an instruction, an output format, a notation, and inputs.

#### A.1.2  Details

The models, inputs, and outputs of each step are shown in Tab.10.

### A.2  Dataset

**Med.rand** consists of randomly chosen questions from PubMedQA. **Med.sim**, on the other hand, is a selection from PubMedQA based on the similarity to the paragraph: *"A medical history of arterial hypertension was associated with lower MMSE scores and a higher prevalence of dementia and cognitive decline at baseline. However, intact cognition through the observation period was linked to higher baseline SBP."* The similarity is determined by comparing the embeddings of this paragraph with those of all other paragraphs in PubMedQA.

### A.3  Human Baseline

#### A.3.1  Human Baseline Description

**Annotator Selection**  We selected three annotators with expertise in the relevant fields to ensure the quality of the questions and annotations. All annotators had prior experience in data annotation and a good understanding of the subject matter.

**Question Formulation**  The annotators were instructed to manually formulate questions based on the original text provided in the datasets. They were asked to create questions that would test the comprehension and response-generation capabilities of the models.

**Annotation Process**  The annotators annotated the model responses while having access to the original text. This approach allowed them to assess the accuracy of the model's answers in the context of the given information.

#### A.3.2  Annotation Scoring

**Scoring Criteria**  Annotations were scored on a binary scale: 0 for incorrect answers and 1 for correct answers. An answer was considered correct if it accurately addressed the question based on the information provided in the original text.

**Scores**  The average scores for each dataset and model are presented in Tab.9:

**Inter-Annotator Agreement**  To ensure the reliability of the annotations, we calculated the inter-annotator agreement using the Fleiss' kappa coefficient. The kappa value was found to be 0.72, indicating substantial agreement among the annotators.

| Model | Score 1 | Score 2 | Score 3 | Score 4 | Score 5 |
|---|---|---|---|---|---|
| GPT-4 | 0.9566 | 0.9479 | 0.8269 | 0.9190 | 0.9230 |
| GPT-3.5 | 0.9245 | 0.9056 | 0.5094 | 0.7735 | 0.8113 |
| Llama2 | 0.8219 | 0.8062 | 0.4002 | 0.6295 | 0.6352 |
| MedLlama | 0.6367 | 0.6241 | - | - | - |

Table 9: Human scores for different models



Figure 6: Prompt for question generation.



Figure 7: Prompt for question answering.

**Prompting of Judging**

**Instruction**

You are now an interviewer and you need to evaluate the candidate's understanding of the domain: [DOMAIN].
Your current task is to judge whether the answer from the candidate is correct or related according to the reference text.
There are 3 types:
"Unrelated" means "The candidate's answer is completely irrelevant to the question and the reference text, and does not provide any relevant information."
"Correct" means "The candidate's answer partially or completely aligns with the reference text. The candidate's answer may contain some omissions, incompleteness or even some details, facts or dimensions that are not mentioned in the reference text, but overall demonstrates the candidate's understanding of the knowledge and their ability to accurately apply and explain relevant concepts, facts, or principles."
"Incorrect" means "The candidate admits that he/she does not know the answer" or "The candidate's answer contradicts, is completely incorrect compared to the reference text. It may contain factual errors, misunderstandings, or incorrect inferences that are inconsistent with the reference text. Pay attention that never make the mistake of thinking that facts you don't know are incorrect! "
The candidate's answer may provide some details, facts or dimensions that are not mentioned in the reference text: You can just assume these details or facts are correct. Because the reference cannot not provide exact factual support to prove that these details or facts are indeed incorrect. However, if all dimensions or details provided by candidate are not mentioned in the reference text, the output type should be "unrelated".

**Output Format**

The output should be formatted as
{
          "type": "unrelated, incorrect or correct",
          "reason": "the reason why you give the type, if the type is incorrect, please point out the exact error from candidate\'s answer and the exact correct answer from the reference text."
}

**Notation**

NOTE:
1. Please output according to the output format in one line.
2. The reason should be limited in 100 words.

**Inputs**

The answer from the candidate is:          [Answer]
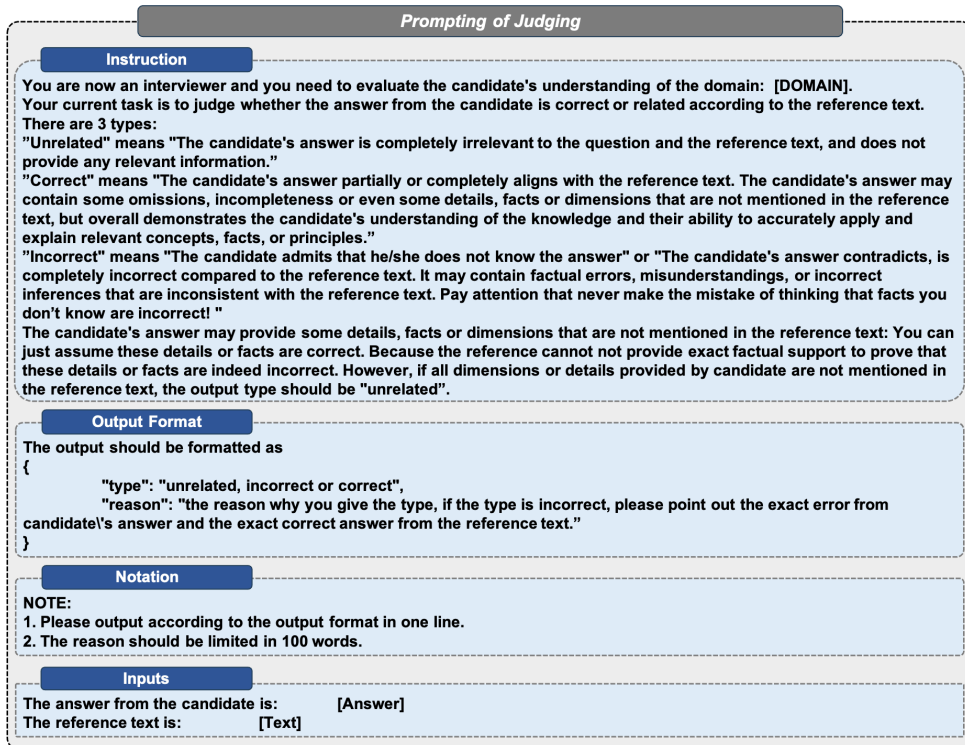The reference text is:          [Text]

Figure 8: Prompt for response judgment.

## A.4 Other Factors in RECKON

**The more knowledge units reduce in each chatting round, the better the LLM's ability to answer relevant knowledge questions.** Fig.9 shows the performance of LLM across multiple iterations of RECKON. The slope of these curves, denoted as the difference between the count of knowledge units in the $(n + 1)$th round and the $n$th round, represents the amount of knowledge units that are reduced in each round. More powerful models like GPT4-turbo tend to answer more questions in each chatting round, especially at the outset and particularly on challenging datasets. While LLaMA2-Chat consistently exhibits slower speeds compared to GPT4-turbo and GPT3.5-turbo, which indicates the low capacity of LLaMA2-Chat for question answering compared to the other models.

**In RECKON, Models with stronger capacity incur fewer token costs during the evaluation process in our method.** Costs of GPT4-turbo and LLaMA2-chat in Tab.4 reflect that the cost of evaluating GPT4 on RECKON is 26% lower than that of LLaMA. This cost advantage shows in the stage of both Question Generstion and Response Judgement. As shown in Fig.9, the amount of knowledge units that are reduced in each round is different between models, and this also causes the difference in cost between models. Specifically, GPT4-turbo has a significant reduction in the number of knowledge units in each round, which decreases the repetition of questions and judgments, thereby reducing the cost of API calls.
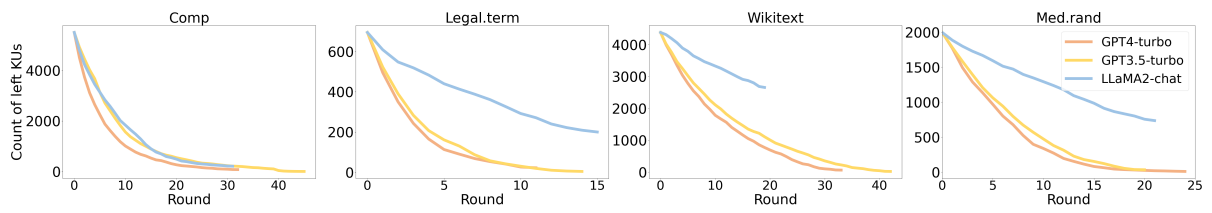
Figure 9: The number of unjudged knowledge points remaining after each iteration of different models on different data sets. The iteration ends when the representation on RECKON converges.

Table 10: The models, inputs, and outputs of each step in RECKON

| Step | Model | Input Example | Output Example |
|---|---|---|---|
| Preparing Knowledge Unit Candidates | gpt3.5-turbo-0125 | **"text"**: "lastaired = ' ' Sa Sandaling Kailangan Mo Ako ' ' ( lit . ' ' The Moment You Need Me ' ' ) is a Philippine drama series aired on ABS-CBN in 1998-2001 . It starred some of the Philippines ' renowned veteran actors and actresses combined with promising young artists under the direction of the reputable Filipino TV and movie director , Olivia M. Lamasan . This shows the intricacies of love , friendship , betrayal and forgiveness. # Cast # *Marvin Agustin as Ruben *Kristine Hermosa as Agnes *Piolo Pascual as Raffy *Giselle Toengi as Stella *John Lloyd Cruz as Daniel *Kaye Abad as Eloisa *Dante Rivero *Hilda Koronel *Ronaldo Valdez *Tessie Tomas *Tracy Vergel *William Lorenzo # See also # *List of programs broadcast by ABS-CBN *List of programs aired by ABS-CBN *List of telenovelas of ABS-CBN" | **"type"**: "Factual Knowledge", **"keyword"**: "Sa Sandaling Kailangan Mo Ako", **"description"**: "A Philippine drama series aired on ABS-CBN from 1998-2001, directed by Olivia M. Lamasan." |
| Clustering | gpt4-turbo-0409 | **"knowledge_units"**: ["bundle to pass data: Passing required data to subfolders using bundle.", "Types of storage in Android: Shared Preferences, Internal storage, External storage, SQLite database, Network connection", "Service: A background component for running music or processing network transactions", "Fragment: A part of an Activity to display multiple screens.", "Sleep mode for CPU: CPU is in a sleep state and not accepting commands except for wireless and alarms.", ...(and 729 more)] | **"center"**: "Web Services: Open standards for cross-platform communication between applications", **"knowledge_units"**: ["Web services: Client and server apps communicate with HTTP protocol for interoperability.", "SOAP and REST: SOAP for data transmission, REST for handling media components.", "REST and Routing in Rails: REST represents representational state transfer, uses HTTP methods to represent actions", "API and Web Services differences: APIs may not require network access, while Web Services always do", "SOAP: XML-based protocol to exchange information between computers.", "REST API: A set of functionalities to perform requests and responses over HTTP.", "REST features: REST is stateless, uses GET method to access resources, and differs from web services.", ...(and 64 more)] ...(more centers and knowledge_units) |

| | | | |
|---|---|---|---|
| Question Genera-tion | gpt4-turbo-0409 | **"center"**: "Exception handling in various programming languages", **"knowledge_units"**: ["Exception handling: Mechanism used to handle runtime errors and maintain normal flow of program.", "Exception Handling: Process of catching and handling exceptions in a program", "Exception Handling in Java: How catch block changes variable value and how exception handling works", "Exception handling in Django: Using exception classes in Django to handle anomalous events", "Exception Handling: Begin-end block with rescue clause specify exception type to handle.", "Exception handling: A mechanism that handles exceptional situations and prevents abnormal program termination.", "Quit command: A command to close the connection.", "rmdir command: Removes the directory specified by the filename.", "Python exception handling: Using try-except blocks to handle errors, with available clauses try-except-finally and try-except-else.", "Exception handling: A mechanism to deal with runtime errors, maintaining the normal program flow.", "Exception handling in Java: Handling and propagation of exceptions in Java programs", "System class: Access system resources, standard input and output.", "gunzip command: a command used to decompress gzip files", "Exception handling in Django: Django uses its own exception classes and supports all core Python exceptions.", "Exception Handling: Using begin-end and rescue clause to handle exceptions", "Debugging: The process of finding and fixing errors in code.", "'type' and 'typeAndWait' command: Enter keyboard key or select value, wait for page reload", ...(and 10 more)] | **"question"**: "In the context of various programming languages, explain how exception handling works in Java, including how a catch block can change a variable's value, the use of try-except blocks in Python with their clauses, and how Django utilizes its own exception classes to handle errors. Additionally, contrast these mechanisms with the general process of catching and managing exceptions to maintain the normal flow of a program. Then tell me what else you know about these.", "center": "Exception handling in various programming languages" |

| | | | |
|---|---|---|---|
| Get Response | model to be evaluated | **"question"**: "In Java object-oriented programming, how does the concept of inheritance differ from implementing an interface, and can you explain the use of the 'javap' command in relation to these concepts? Also, can you describe the role of JavaBeans in encapsulating state and behavior within an object, the differences between singleton and prototype bean scopes in Spring, and how thread creation can be achieved? Please include an explanation of how a thread can be interrupted and the implications of calling 'join()' on a thread, along with the considerations for thread-safety when dealing with singleton patterns. Then tell me what else you know about these." | **"answer"**: "In Java, inheritance is a mechanism where one class acquires the properties (methods and fields) of another class. The child class extends the parent class, and it can use all the public and protected methods and fields of the parent. On the other hand, implementing an interface provides a way for a class to become more formal about the behavior it promises to provide. It specifies a set of methods that the class must implement. Unlike inheritance, a class can implement multiple interfaces, but it can only extend one class. The 'javap' command in Java is used to print information about the classes and interfaces in a given Java binary file. It displays the package, protected, and public fields and methods of the classes available in the specified file. JavaBeans are reusable software components that adhere to specific naming conventions for properties, methods, and events. They provide a way to encapsulate state and behavior within an object. Singleton and prototype bean scopes are used in Spring to control the lifecycle of objects. Singleton beans are created only once per container, while prototype beans are created every time they are requested. Thread creation in Java can be achieved by either extending the Thread class or implementing the Runnable interface. Interrupting a thread is done by calling the interrupt() method of the thread object. Calling join() on a thread waits for it to complete before moving on to the next task. Thread-safety considerations are important when dealing with singleton patterns because multiple threads may access the same instance concurrently. In addition to these topics, I have knowledge of other Java concepts such as generics, collections, exception handling, JDBC, JPA, Hibernate, Servlets, JSP, RESTful web services, Spring Framework, and more." |

| Judging | gpt3.5-turbo-0125 | **"knowledge_unit"**: "Java Bean: A reusable software component that encapsulates many objects into one object for easy maintenance.", **"answer"**: "In Java, inheritance is a mechanism where one class acquires the properties (methods and fields) of another class. The child class extends the parent class, and it can use all the public and protected methods and fields of the parent. On the other hand, implementing an interface provides a way for a class to become more formal about the behavior it promises to provide. It specifies a set of methods that the class must implement. Unlike inheritance, a class can implement multiple interfaces, but it can only extend one class. The 'javap' command in Java is used to print information about the classes and interfaces in a given Java binary file. It displays the package, protected, and public fields and methods of the classes available in the specified file. ...(more)" | **"knowledge_unit related content"**: "JavaBeans are reusable software components that adhere to specific naming conventions for properties, methods, and events. They provide a way to encapsulate state and behavior within an object. Singleton and prototype bean scopes are used in Spring to control the lifecycle of objects. Singleton beans are created only once per container, while prototype beans are created every time they are requested." |
| | gpt4-turbo-0409 | **"knowledge_unit related content"**: "JavaBeans are reusable software components that adhere to specific naming conventions for properties, methods, and events. They provide a way to encapsulate state and behavior within an object. Singleton and prototype bean scopes are used in Spring to control the lifecycle of objects. Singleton beans are created only once per container, while prototype beans are created every time they are requested.", "text": "Singleton scope should be used together with EJB stateless session bean, and prototype scope should be used together with EJB stateful session bean." | **"type"**: "incorrect", **"reason"**: "The candidate incorrectly associates singleton and prototype bean scopes with Spring, while the reference text specifically mentions their use with EJB (Enterprise JavaBeans) stateless and stateful session beans, not Spring." |