# ADDRESSING THE SUSTAINABLE AI TRILEMMA: A CASE STUDY ON LLM AGENTS AND RAG

**Hui Wu**
Dept. of Engineering
University of Exeter
Exeter, UK
EX4 4QF
hw865@exeter.ac.uk

**Xiaoyang Wang**
Dept. of Computer Science
University of Exeter
Exeter, UK
EX4 4QF
x.wang7@exeter.ac.uk

**Zhong Fan**
Dept. of Engineering
University of Exeter
Exeter, UK
EX4 4QF
z.fan@exeter.ac.uk

January 15, 2025

## ABSTRACT

Large language models (LLMs) have demonstrated significant capabilities, but their widespread deployment and more advanced applications raise critical sustainability challenges, particularly in inference energy consumption. We propose the concept of the Sustainable AI Trilemma, highlighting the tensions between AI capability, digital equity, and environmental sustainability. Through a systematic case study of LLM agents and retrieval-augmented generation (RAG), we analyze the energy costs embedded in memory module designs and introduce novel metrics to quantify the trade-offs between energy consumption and system performance. Our experimental results reveal significant energy inefficiencies in current memory-augmented frameworks and demonstrate that resource-constrained environments face disproportionate efficiency penalties. Our findings challenge the prevailing LLM-centric paradigm in agent design and provide practical insights for developing more sustainable AI systems. Our code is available online[1].

## 1 Introduction

Large language models (LLMs) have played an important role in artificial intelligence, significantly advancing natural language understanding and generative AI. These models are now powering applications across healthcare, education, and business, demonstrating their growing practical impact in everyday services. However, this rapid development comes with significant challenges, particularly in terms of energy use and its impact on the environment and society. While the energy consumption of LLM training has received much attention [1] , recent evidence [2] [3] [4] suggests that the inference phase, particularly in large-scale deployments, has emerged as the dominant contributor to the overall energy footprint, highlighting a critical and growing concern in the sustainability of LLM-powered services.

In response to these challenges, we propose the concept of the **Sustainable AI Trilemma**, which draws parallels to the sustainability trilemma balancing Economy, Society, and Environment [5]. In the AI context, this trilemma highlights the trade-offs among **capability, digital equity, and environmental sustainability**. Current LLM deployments heavily prioritize capability—pushing for greater model performance and broader applications—often at the expense of equity and sustainability. This imbalance is made worse by the increasing use of **LLM-based agents**, and the scaling of **multi-agents**, which applies LLMs in a frequent manner for long-term operations.

LLM-based agents, capable of executing complex, multi-step tasks, are increasingly integrated across domains such as enterprise automation, scientific research, and personalized assistance [6]. This paradigm amplifies the energy demands of LLM inference, raising urgent questions about the sustainable scalability of such AI services. Without deliberate attention to energy efficiency, the widespread use of these agents risks accelerating resource inequities [7] and undermining global sustainability goals [8].

---

[1]https://github.com/huiwxing/llmagent_trilemma

Despite these challenges, most research on LLM agents focuses on improving performance, with little consideration of sustainability or equity metrics. As a first step toward addressing the trilemma, we use LLM agents as a case study to better understand these challenges. This paper aims to address this gap with the following contributions:

**Proposing the Sustainable AI Trilemma and Highlighting its Urgency:** We articulate the risks posed by an overemphasis on productivity in LLM agent deployments, which exacerbates the trilemma by neglecting equity and environmental sustainability.

**Unified Analysis of Energy Cost in LLM Agents:** Through systematic analysis and a case study of the memory module, we reveal the energy costs embedded in LLM agent and RAG designs. Our findings emphasize the trade-offs between enhanced functionality and sustainability, demonstrating that resource-intensive designs often yield diminishing performance returns and identifying opportunities for more sustainable implementations.

**New Metrics for Energy and Performance Trade-offs:** We introduce novel metrics to evaluate the trade-offs between energy cost and performance in LLM agents. These metrics provide a foundation for optimizing the design and operation of LLM agents to achieve more energy efficiency without compromising effectiveness.

**Challenging the Paradigm of LLM-Centric Autonomy:** A growing trend in agent design involves offloading all subtasks to the LLMs, purportedly to enhance autonomy. From the perspectives of environmental impact and cost-efficiency, we challenge this prevailing trend, demonstrating through empirical evidence that this approach often results in unnecessary computational overhead without proportional benefits in agent autonomy or performance, relying excessively on LLMs undermines the broader objectives of sustainable AI.

The remainder of this paper is organized as follows: Section 2 introduces the Sustainable AI Trilemma and its four constituent dilemmas, along with background on LLM agents and LLM energy consumption. Section 3 presents a case study on memory modules in LLM-based agents, examining their impact on energy efficiency and system complexity, and proposes metrics for evaluating energy-performance trade-offs. Section 4 presents our experimental results and analysis. Finally, Section 5 discusses implications and concludes with insights for future sustainable AI development.

## 2   Background

### 2.1   Sustainable AI Trilemma

The concept of a trilemma, which highlights the difficulty of balancing competing priorities, has been explored in various contexts, including sociology and ethics. Recent sociological studies have applied this trilemma framework to AI [9], and others have emphasized the sustainability of AI [10][11]. The AI trilemma posits that AI development must balance three fundamental goals: AI capability or performance, sustainability, and equity, where optimizing for any two often comes at the expense of the third. However, with the rapid rise of LLMs and their pervasive influence, we argue that this trilemma must be redefined and addressed from a technical perspective. The Sustainable AI Trilemma emerges from a series of intertwined technical dilemmas that reflect conflicting priorities in the design, deployment, and long-term operation of AI systems. Below, we identify and analyze four core dilemmas, each highlighting a critical technical challenge.

**Evaluation Dilemma: Awareness vs. Measurement**   While the environmental impacts of LLMs and other AI systems have become a recognized concern, the lack of accurate and comprehensive assessment tools hinders progress toward sustainability. Existing tools, such as MLCO2 [12] and Green Algorithms [13], are designed for traditional machine learning and earlier AI algorithms. These tools fail to account for the unique challenges posed by LLMs and their applications, such as memory-intensive inference processes, highly variable usage patterns, and the lifecycle emissions from large-scale deployment. This technical gap leaves practitioners and policymakers uncertain about the true scale of AI's environmental footprint, creating a paradox where actionable insights are hindered despite widespread recognition of the problem.

**Design Dilemma: Autonomy vs. Efficiency**   The pursuit of highly autonomous AI systems [6], exemplified by the rise of single- and multi-agent frameworks, drives the demand for increasingly sophisticated architectures. While these systems promise greater flexibility and functionality, their designs often introduce redundancy and inefficiencies [14], such as poorly optimized workflows leading to excessive and unnecessary model calls, frequent over-reliance on LLMs for subtasks that could be addressed more efficiently with specialized modules, or maintaining multiple concurrent LLM instances for similar tasks across different services. This tradeoff between expanding autonomy and maintaining energy-efficient operations adds strain to system-level resource management, resulting in diminishing returns in capability.
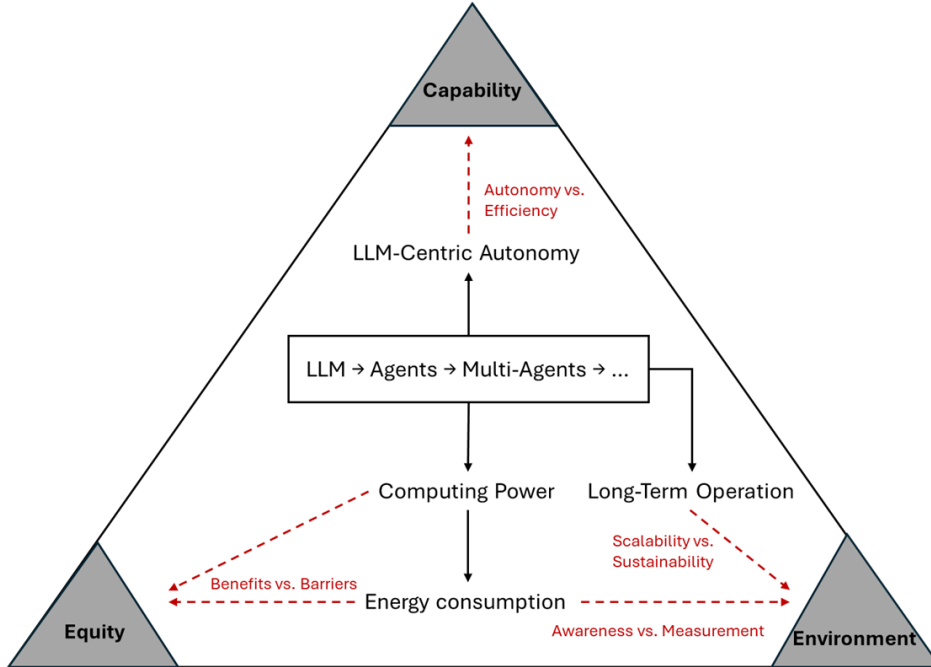
Figure 1: The Sustainable AI Trilemma.

**Operational Dilemma: Scalability vs. Sustainability**    Scaling AI services to meet increasing demand exacerbates the tension between scalability and sustainability. Expanding the number of deployed LLM-based agents requires maintaining consistent Quality of Service (QoS) under growing workloads. However, the infrastructure required to ensure QoS—such as low-latency inference, real-time responsiveness, aggressive caching and pre-warming strategies [15]—inevitably leads to heightened energy consumption and operational costs. At scale, these architectural inefficiencies can be further amplified. While traditional AI systems benefit from MLOps and AIOps [16] to streamline operations, these tools face significant challenges when applied to the complexity of LLM applications. The memory-intensive processes, cross-agent interactions, and continuous learning demands of LLM-based systems strain the existing operational workflows, raising questions about their scalability and long-term feasibility in resource-constrained environments.

**Access Dilemma: Benefits vs. Barriers**    Resource-intensive AI applications disproportionately benefit well-resourced organizations and regions while marginalizing underprivileged communities. Proprietary LLMs often rely on expensive API-based access models, imposing high financial costs on users. Open-source LLMs, while enabling local deployment, demand highly efficient hardware to mitigate their energy-intensive operations, which are often unavailable in resource-limited regions. These twin pressures—economic costs and environmental impacts—place an outsized burden on communities with limited access to affordable and reliable computing infrastructure. The resulting disparities exacerbate the digital divide [7], effectively excluding these communities from accessing and benefiting from advanced AI services that could otherwise drive meaningful societal and economic change [17][18].

These four dilemmas collectively form the foundation of the **Sustainable AI Trilemma**: the need to balance AI capability, digital equity, and environment. To address this trilemma, we believe a paradigm shift in designing, deploying, and evaluating AI systems is required.

## 2.2   LLM-based Agents

The emergence of LLM-based agents represents a transformative shift in the application of large language models. Unlike traditional applications that rely on single-turn interactions, LLM-based agents are designed to perform multi-turn dialogues and execute complex, multi-step tasks autonomously [19]. These agents are increasingly deployed across diverse sectors, from intelligent customer support systems to advanced research assistants, showcasing their potential to revolutionize workflows and decision-making processes. At their core, LLM-based agents integrate multiple functional components to achieve greater autonomy and adaptability. Key components typically include:

- **Planning Module**: Breaks down high-level tasks into manageable sub-tasks, enabling structured and goal-oriented operations [20].
- **Tool Use**: Incorporates external APIs, databases, or specialized software tools to enhance problem-solving capabilities [21].
- **Memory Module**: Allows agents to retain and retrieve past information, supporting contextual understanding and long-term consistency.

Among these, the memory module is particularly important. While traditional LLM applications depend solely on the pretrained knowledge embedded within the model and the immediate context provided by user queries, memory-enabled LLM agents or Retrieval-Augmented Generation (RAG) [22] extend this paradigm by incorporating both internal memory (e.g., conversation histories) and external memory (e.g., indexed document repositories) [23]. This memory capability facilitates advanced applications by mitigating LLM hallucination [24], providing improved performance through enriched context and persistent knowledge access.

### 2.3 Energy Cost and Environmental Impacts of LLM Inference

The growing energy consumption of AI systems, particularly LLMs, has become a significant concern in the field and inference-phase energy consumption deserves particular attention, especially in terms of carbon footprints [25]. [3] suggests that the multiplicative effect due to widespread AI deployment remains a crucial consideration. Studies in [2] have shown that for services like ChatGPT, inference-related carbon emissions in one year amount to 25 times that of GPT-3's training. This trend has driven the development of specialized measurement tools and analysis. Green Algorithms [13] provides an online platform for estimating computational carbon footprints, while LLMCarbon [8] specifically addresses LLM scenarios, offering end-to-end carbon emission predictions for both dense and mixture-of-experts (MoE) architectures.

Empirical studies have extensively investigated LLM inference energy consumption. [26] have conducted large-scale testing of LLaMA models across different GPU platforms (V100 and A100), including scenarios with up to 32-GPU model sharding. Through comprehensive characterization studies of mainstream LLMs, Wilkins et al. [27] have developed high-accuracy ($R^2 > 0.96$) models for energy consumption and runtime prediction. Further research [28] has explored how factors such as model size, layer count, parallel attention mechanisms affect energy consumption, along with optimization techniques like batch processing and quantization. In code generation, Vartziotis et al. [29] have evaluated the sustainability awareness of commercial AI models like GitHub Copilot, ChatGPT-3, and Amazon CodeWhisperer. What's more, the optimization of AI energy consumption requires careful consideration of trade-offs - while some AI models aim to reduce carbon emissions, their training and deployment also generate environmental impacts, yielding positive net benefits only under appropriate scenarios [30].

## 3 Computational and Energy Costs in LLM-based Agents: A Case Study on Memory Modules

The memory module is a transformative feature in LLM agents. However, the incorporation of memory introduces notable challenges in energy efficiency and design complexity. In this section, we take the memory module of LLM-based agents as an example to explore potential inefficiencies in resource utilization within their design and operations.

### 3.1 Problem Definition

To understand its impact, we first define the simplest case of an LLM agent without memory before formalizing the enhanced capabilities of a memory-augmented agent.

**Without Memory** In its simplest configuration, an LLM agent operates as a standalone inference system. Given a user query $Q$, the LLM generates a response $A$:

$$LLM(Q) = A \tag{1}$$

This formulation assumes no access to external data or stored context. The model relies entirely on its pretrained knowledge and immediate input $Q$. While computationally efficient, this approach is limited by the LLM's inherent knowledge boundaries, which cannot be updated without retraining. When $Q$ references information outside the model's scope, the output $A$ is prone to inaccuracies or hallucination, failing to meet the requirements of knowledge-intensive or context-dependent applications.

**With Memory** When a memory module is incorporated, the memory $M$ is typically categorized into internal memory and external memory, i.e., $M = \{M_{internal}, M_{external}\}$. Internal memory comprises historical information collected during the agent's operation. For example, an agent solving iterative tasks may store user inputs $i$ and its own responses $o$ from prior interactions at time $t$: $M_{internal} = \{(i_1, o_1), (i_2, o_2), \ldots, (i_{t-1}, o_{t-1})\}$. External memory refers to data retrieved from domain-specific repositories, databases, documents, or other external sources: $M_{external} = \{D_1, D_2, \ldots, D_n\}$, where these external resources $D$ are indexed and queried as needed. The behavior of the memory-enabled agent can be described as:

$$LLM(Q, M_q) = A \tag{2}$$

Here, $M_q \in M$ represents the relevant subset of the memory $M$, retrieved based on the query $Q$. We assume that for every query there is a ground truth answer $A_G$, which can be generated based on their ground truth memories $M_G$.

**Energy Model** Prior work [27] has established an analytical model for LLM inference energy consumption that accounts for both input and output text lengths, as well as their interaction. Specifically, the energy consumption $E$ of a primary input-output of a certain LLM $K$ can be expressed as

$$E = e_K(T_{in}, T_{out}) = \alpha_{K,0} * T_{in} + \alpha_{K,1} * T_{out} + \alpha_{K,2} * T_{in} * T_{out} \tag{3}$$

where $T_{in}$ and $T_{out}$ represent the number of input and output tokens respectively, and $\alpha_{K,0}$, $\alpha_{K,1}$, and $\alpha_{K,2}$ are model-specific ($K$) coefficients. Since output token generation is more energy-intensive than input processing [27], $\alpha_{K,1}$ normally exceeds both $\alpha_{K,0}$ and $\alpha_{K,2}$. We will use this model to perform some theoretical simulation analyses in the next section and some metrics calculations in Experiment 4.3.

## 3.2 Cost of Existing Memory Module Design Patterns

The integration of memory modules introduces more challenges. Many existing memory designs completely rely on LLM operations (as shown in Figure 2), prioritizing performance and autonomy improvements without considering computational workloads and resource consumption. Through analyzing different components of memory-augmented frameworks, we identify three critical phases in modern LLM memory architecture: (a) **Memory Formation**, which processes and organizes internal/external input data into storage structures; (b-c) **Memory Reading**, encompassing retrieval necessity detection and query optimization for better retrieval; and (d-f) **Memory utilization**, to generate better answers by optimizing memory usage. This paper aims to uncover the hidden costs of these components and quantify the trade-offs between enhanced task performance and increased energy consumption.

### 3.2.1 Memory Formation

Memory formation involves the transformation of raw data into structured and retrievable representations using various indexing and storage optimization techniques. Vector-based storage and retrieval, which transforms raw data into vector embeddings for efficient similarity-based searches, struggles with scalability and complex query requirements such as comprehensive full-text searches, relation-based graph, and keyword search [31]. Several alternative LLM indexing strategies have been proposed. Keyword or topic-based indexing, implemented in systems like MemoChat [32] and MemoryBank [33], leverages LLMs to extract salient terms from the input data. Summary-based indexing, adopted by systems such as SCM [34], MemGPT [35], and MPC [36], employs LLMs to generate concise representations of memory segments. A more sophisticated approach is to build structured representations, such as triplets and knowledge graphs, by extracting entities and relationships through LLM. Systems like Think-in-Memory [37] and RET-LLM [38] exemplify this approach.

While advanced memory representations excel in supporting complex reasoning tasks, they impose significant computational overhead due to multiple LLM processing iterations, raising concerns about scalability and efficiency as memory data volumes grow. This raises important questions about the trade-off between costs and performance benefits. For instance, when most of indexed memories see limited or no retrieval, the substantial energy investment in memory formation processes becomes difficult to justify against the minimal practical benefits. Such scenarios highlight the critical challenge of aligning indexing strategies with specific application requirements to ensure optimal resource utilization.

### 3.2.2 Memory Reading

Through the preprocessing process of Memory Formation, the LLM agent is now equipped with a retrievable memory/knowledge base. Memory Reading, on the other hand, is the process of receiving a query and then getting a portion of information from the base that matches the query.
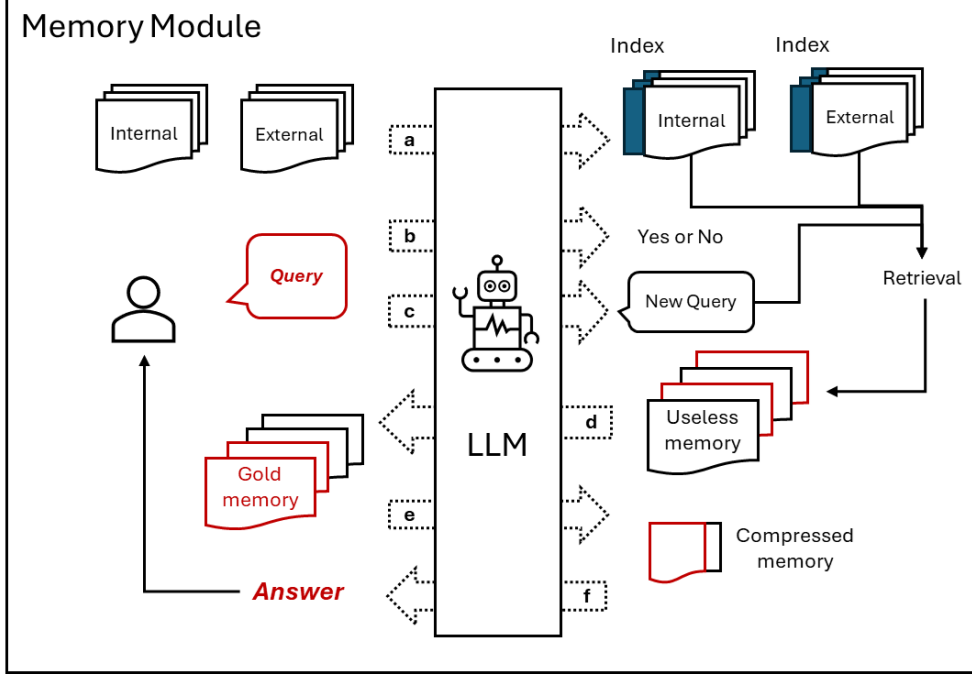
Figure 2: Design of Memory Modules in LLM Agents or RAG. The figure contains the operations that require LLM inference: (a) Memory Formation, (b) Retrieval Necessity Detection, (c) Query Optimization, (d) Reranking, (e) Compression, (f) Generation. The red markers represent data that is necessary for a QA process to generate a correct response, while black indicates non-necessary data.

**Retrieval Necessity Detection**   While RAG improves response accuracy and reduces hallucinations, persistent indiscriminate retrieval can cause unnecessary overhead. Retrieval necessity detection optimizes efficiency by using the LLM's inherent knowledge for simple tasks and enabling selective memory access when necessary. SCM [34] and Ren et al. [39] design prompt like Self-Asking prompt and Judgment prompt to assess memory activation needs. MemGPT [35] treats LLM as an operating system and empowers the LLM processor to autonomously decide, update, and search its memory based on the current context. Rowen [40] introduces a multilingual semantic-aware detection module that evaluates cross-language response consistency and triggers retrieval augmentation when factual inconsistencies are detected. SKR [41] uses both direct prompting and training a smaller classifier, to detect whether the model has "self-knowledge".

Building on these studies, we distill the Retrieval Necessity Detection process to its fundamental form, as outlined below in terms of the LLM inference operations (To simplify the analysis, system prompt templates in the input are temporarily excluded. The following analysis remains consistent.):

- **Input**: Query (token number = $T_q$).
- **Output**: A boolean or numeric score (indicating whether retrieval is needed, token number = $C$, where $C$ is a small constant).
- **Energy**: $E_1 = e_K(T_q, C) = \alpha_{K,0} * T_q + \alpha_{K,1} * C + \alpha_{K,2} * T_q * C$

**Query Optimization**   Query optimization enhances retrieval accuracy by enriching raw or under-contextualized queries with semantic information through query expansion or transformation. Query expansion generates additional queries derived from the original input like decomposing complex questions into sub-questions, or employing prompt engineering techniques to create multiple contextually enriched variations, such as the Least-to-most prompting [42], Chain-of-Thought [43], Self-Ask [44], and Multi Query Retrieval [45]. Query transformation focuses on rewriting or reinterpreting the original query to better align with retrieval objectives, such as Hypothetical Document Embedding (HyDE) [46], Step-back prompting [47], Rewrite-Retrieve-Read [48]. Both approaches increase the semantic richness and correctness of the query, enabling retrieval systems to access a broader and more relevant information space. However, it is important to note that these enhancement methods typically not only increase the token count of the original query but also represent an additional computational cost in themselves.

The basic form of query optimization is as follows:

- **Input**: Query (token number = $T_q$).
- **Output**: New query/queries, or hypothetical documents (token number = $T_{qn}$, which is usually larger than $T_q$).
- **Energy**: $E_2 = e_K(T_q, T_{qn}) = \alpha_{K,0} * T_q + \alpha_{K,1} * T_{qn} + \alpha_{K,2} * T_q * T_{qn}$

**Memory Retrieval**  Memory retrieval is the core step where relevant memory entries are extracts based on the transformed query, typically using similarity metrics (e.g., cosine similarity) between query and document embeddings. Methods include sparse retrieval (e.g., BM25 [49]), dense retrieval (e.g., Sentence-BERT [50]), or hybrid retrieval combining both. Although some approaches like Memochat [32] attempt to use LLMs as retrievers, we argue that such methods are significantly limited by the LLM's context window size and lack scalability in retrieval speed and volume, making them unsuitable for consideration in this study. As a fixed computational cost **independent of LLM inference**, memory retrieval is considered an essential step in the pipeline.

### 3.2.3  Memory Utilization

Memory utilization involves additional processing of retrieved memory, such as re-ranking and compression, to integrate the memory with the query in an appropriate format as input to LLM for generating answers.

**Reranking**  Reranking serves as a refinement step in the RAG process, where retrieved documents are reordered based on their relevance to the query. Recent advancements have increasingly focused on LLM-based reranking. Zhuang et al. [51] demonstrated that LLMs pre-trained on unstructured text, like LLaMA and Falcon, exhibit strong zero-shot ranking capabilities. These methods such as LRL [52], UPR [53], RankGPT [54], PRP [55] leverage prompt engineering with models like GPT-3/4, UL2, and T5 for zero-shot passage reranking, employing various prompt types, including query generation, relevance genenration [56], pairwise ranking, and permutation for ranked lists. Although prompt-based approaches reduce fine-tuning overhead and enhance agent autonomy, they themselves introduce extra computational costs during inference by processing multiple retrieved documents simultaneously. However, similar to Retrieval Necessity Detection, considering that reranker outputs are typically concise (ranked indices or scores, a small constant of output token numbers), these approaches have the potential to be cost-effective with moderate number of documents retrieved.

The basic form of LLM reranker is as follows (if using permutation generation prompting):

- **Input**: New query or queries, Retrieved related memories (token number = $T_{qn} + T_m$).
- **Output**: Ranked index (token number = $K$, where $K$ represents the number of top-K retrieved entries).
- **Energy**: $E_3 = e_K(T_{qn} + T_m, K) = \alpha_{K,0} * (T_{qn} + T_m) + \alpha_{K,1} * K + \alpha_{K,2} * (T_{qn} + T_m) * K$

**Compression**  After retrieval and reranking, relevant documents may still exceed LLM input capacity and introduce noise. While compression strategies can help by extracting essential information, current approaches face significant trade-offs. LLM-based summarization and specialized compression models (e.g. PRCA [57], RECOMP [58]) achieve good performance, but generate significantly more output tokens and consume more energy (owing to the large output coefficient $\alpha_{K,1}$ and the number of output tokens $\beta * T_m$, and the square interaction term $\beta * T_m^2$ ) than simpler steps like retrieval detection or reranking. Alternative entropy-based methods like Selective Context [59] and LLMLingua [60] reduce costs by avoiding full LLM inference, but may compromise compression quality [61]. The overall efficiency gains versus computational overhead remains a critical consideration, especially at scale.

The basic form of LLM compressor is as follows:

- **Input**: Ranked memories, which has the same token as unranked just in a different order (token number = $T_m$).
- **Output**: Compressed memories (token number = $\beta * T_m$, $0 < \beta < 1$, where $\beta$ is the compression ratio).
- **Energy**: $E_4 = e_K(T_m, \beta * T_m) = \alpha_{K,0} * T_m + \alpha_{K,1} * \beta * T_m + \alpha_{K,2} * \beta * T_m^2$

**Generation**  The final step synthesizes the query and compressed memory to produce a coherent answer. The input query and output answer are unavoidable computational costs, making compressed memories the key variable factor. While a larger compression ratio ($\beta$) increases energy costs, it also preserves more information, potentially leading to more accurate and complete answers. The optimal compression ratio thus needs to balance the trade-off between computational efficiency and answer quality based on specific application requirements.

- **Input**: Original query (token number = $T_q$), Compressed memories (token number = $\beta * T_m$, $0 < \beta < 1$) .

- **Output**: Answer (token number = $T_a$).
- **Energy**: $E_5 = e_K(T_q + \beta * T_m, T_a) = \alpha_{K,0} * (T_q + \beta * T_m) + \alpha_{K,1} * T_a + \alpha_{K,2} * (T_q + \beta * T_m) * T_a$

### 3.2.4 Memory Management

Beyond basic memory read and write operations, some research focuses on the long-term maintenance and management of memory in LLM agents. These approaches often emphasize agent autonomy, delegating much of the memory management process to the LLM itself. For instance, Think-in-Memory [37] introduces prompts designed specifically for memory merging and forgetting. These methods also introduce more uncertainties in energy consumption over long-term deployment. Further exploration in our future research to balance autonomy and sustainability in memory management.

## 3.3 Metrics for Energy and Performance Trade-offs

In the evaluation of RAG systems, assessment tasks are typically divided into two phases: **Retrieval** and **Generation**. Corresponding to the evaluation subjects defined in Section 3.1, we consider the following elements: query ($Q$), answer ($A$), ground truth answer ($A_G$), retrieved relevant memories ($M_q$), and ground truth memories ($M_G$). Yu et al. [62] synthesized current research on RAG evaluation and categorized various metrics into five fundamental types: **Relevance**($M_q \leftrightarrow Q$) and **Accuracy**($M_q \leftrightarrow M_G$) in the Retrieval phase, **Relevance**($A \leftrightarrow Q$), **Faithfulness**($A \leftrightarrow M_q$), and **Correctness**($A \leftrightarrow A_G$) in the Generation phase, where $\leftrightarrow$ denotes the alignment or relationship between two elements (e.g., $M_q \leftrightarrow Q$ represents how relevant the retrieved documents are to the query).
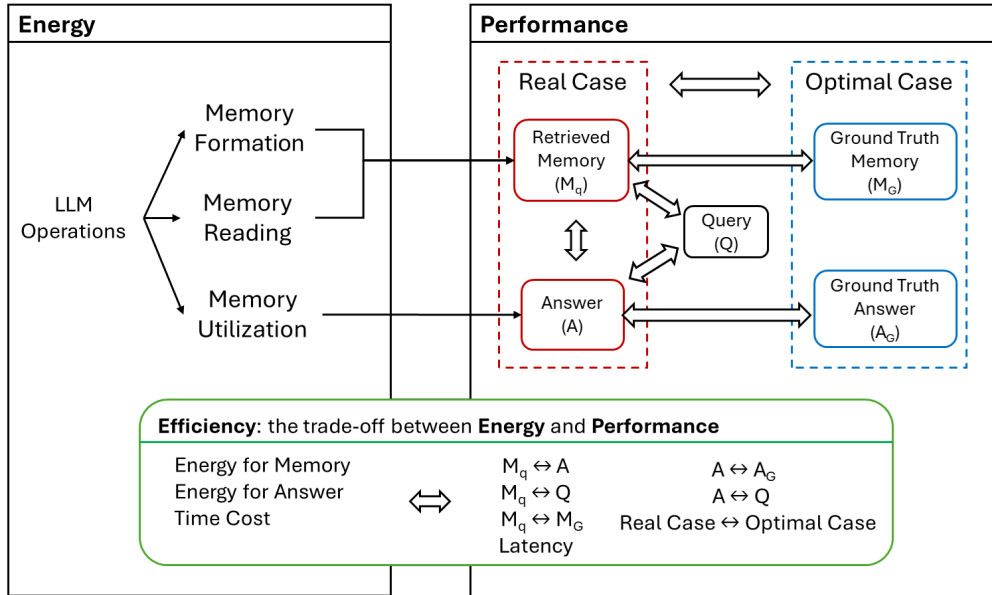


Figure 3: New efficiency metrics to balance energy cost and task performance.

### 3.3.1 Trade-offs Between Memory Formation and Retrieval

The choice of indexing methods during memory formation not only affects system energy consumption but also directly impacts subsequent retrieval efficiency and effectiveness. To comprehensively evaluate different indexing strategies, we propose a multi-dimensional evaluation. For retrieval effectiveness, we employ two key metrics: $Relevance$ measures how well a retrieved memory matches a query and is typically calculated by having the LLM judge each entry's relevance, then dividing the number of relevant entries by the total retrieved. $Accuracy$ is evaluated by comparing retrieved documents to a ground truth set and metrics such as recall and F1 scores can often be used. For retrieval efficiency, we use the average retrieval time **per memory token** $T_{retr}$. For the memory formation process, we focus on energy consumption $E_{token\_m}$ and processing time $T_{cost}$ per memory token.

Based on these metrics, we design two comprehensive evaluation indicators that quantify the energy cost required to achieve a unit improvement in retrieval relevance and accuracy:

**Retrieval Energy-to-Relevance Ratio (RERR)**

$$RERR = \frac{E_{token\_m}}{Relevance} \tag{4}$$

**Energy-to-Accuracy Ratio (EAR)**

$$EAR = \frac{E_{token\_m}}{Accuracy} \tag{5}$$

### 3.3.2 Trade-offs between Memory Reading, Utilization and Generation

**a. For Correct Generations**  For correct generations of the query $Q$ (token number = $T_q$), we propose an ideal case scenario that achieves accurate answers through minimal necessary data (as shown in Figure 2) and operations: in this optimal case, the system would use only the original query with vector embeddings, employing sparse or dense retrieval to extract precisely all relevant ground truth memories $M_G$ (token number = $T_{mg}$) and no irrelevant memories, which are then combined with the original query to directly generate the ground truth answer $A_G$ (token number = $T_{ag}$) through the LLM, without requiring any additional enhancement operations (such as query optimization, reranking or compression).

In the optimal case, the essential energy consumption comprises two components: retrieval cost ($E_{retr}$) and response generation cost $E_g = e_K(T_q + T_{mg}, T_{ag})$. Thus, the optimal total energy consumption can be expressed as:

$$E_{optimal} = E_{retr} + E_g \tag{6}$$

However, in real cases, to ensure generation quality, various additional memory reading and utilization operations in 3.2.2 and 3.2.3 are often necessary. We define $E_{real}$ as the cumulative energy consumption across all stages of the generation process, for example, when using all operations, $E_{real} = E_{retr} + E_1 + E_2 + E_3 + E_4 + E_5$.

To quantify the energy efficiency of the generation process, we propose:

**Generation Energy Optimality Ratio (GEOR)**

$$GEOR = \frac{E_{optimal}}{E_{real}} \tag{7}$$

This ratio serves as an indicator of how close the system's actual energy consumption is to the theoretical minimum. A GEOR value closer to 100% indicates higher energy efficiency, suggesting that the system's energy consumption approaches the ideal case, while lower values indicate potential areas for energy optimization in the generation pipeline.

**b. For All Generations**  The effectiveness and latency of response generation are significantly influenced by memory reading and utilization methods. For generation effectiveness, we employ three key metrics: $Relevance$ measures how well the generated response aligns with the intent and content of the initial query $Q$, typically calculated as a normalized score assigned by an LLM judge based on predefined assessment criteria. $Faithfulness$ evaluates the consistency between the generated response and the retrieved memories $M_q$, calculated as the ratio of main points in the answer attributable to the context to the total main points in the answer. $Correctness$ assesses the accuracy of the generated response against ground truth $A_G$, measured using standard NLP metrics such as ROUGE, BLEU, or BERTScore. For generation efficiency, we measure the average $latency$ per response token and energy consumption per token $E_{token\_g}$ which is the division of $E_{real}$ by **response token** number. Similarly, there are three energy efficiency ratios:

**Generation Energy-to-Relevance Ratio (GERR)**

$$GRER = \frac{E_{token\_g}}{Relevance} \tag{8}$$

**Energy-to-Faithfulness Ratio (EFR)**

$$FER = \frac{E_{token\_g}}{Faithfulness} \tag{9}$$

**Energy-to-Correctness Ratio (ECR)**

$$CER = \frac{E_{token\_g}}{Correctness} \tag{10}$$

These ratios quantify the energy cost required to achieve unit improvement in each generation effectiveness dimension. Lower values indicate better energy efficiency in achieving the respective effectiveness goals.

## 4 Experiments

### 4.1 Experiment settings

**Datasets.** For Internal Memory, we treat the task as that of a daily assistant agent and utilize the LoCoMo [63] dataset, which comprises very long-term dialogues. For External Memory, we use HotpotQA [64] and MuSiQue [65], which are well-known datasets for evaluating multi-hop and complex reasoning abilities in retrieval-augmented tasks. Additionally, we incorporate the AI ArXiv dataset [66], with 100 pairs of queries and answers on topics of RAG extracted by Eibich et al. [67]. We follow the same node splitting configuration as in Eibich's experiments. For all datasets we sampled some data and the characteristics of these sub-datasets are shown in Table 1.

Table 1: Comparison of Dataset Characteristics

| Dataset | Memory Node | | Query | | | Answer | Scenario |
|---|---|---|---|---|---|---|---|
| | Number | Size* | Types | Difficulty | Length | Length | |
| LoCoMo | 419 | 54.6 | 5 | Easy | Short | Short | Daily Life |
| HotpotQA | 992 | 128.5 | 1 | Hard | Short | Short | Open-Domain |
| MuSiQue | 2000 | 111.7 | 1 | Hard | Short | Short | Open-Domain |
| AI ArXiv | 51270 | 27.5 | 1 | Medium | Long | Long | Academic |

*The size of memory node is the average token number.

**Metrics tools.** We evaluate Retrieval-Relevance, Generation-Relevance and Faithfulness using UpTrain's [68] LLM-based scorer with the GPT-4 API. For Retrieval-Accuracy we measure the Mean Reciprocal Rank (MRR) and Recall based on the ground truth provided by the datasets, while Generation-Correctness is measured using Recall-Oriented Understudy for Gisting Evaluation (ROUGE) [69]. All score ranges from 0 to 100.

**Hardware.** To evaluate LLM performance across different resource scenarios, we configured two hardware setups. The **resource-constrained environment** utilizes a single NVIDIA GeForce RTX 4080 GPU(16GB, 320W TDP) and an Intel i7-13700K processor (16-core 24 threads, 125W TDP). For the **resource-abundant environment**, we employed a cluster of eight NVIDIA A100 GPUs (80GB, 400W TDP) and two AMD EPYC 7J13 processors (64-Core 128 threads, 280W TDP), with NVIDIA NIM inference accelerating microservices [2], representing advanced high-performance computing capabilities compared with resource-constrained setup. Tasks 1 and 2 are conducted on the resource-constrained setup, while Task 3 leverages the resource-abundant configuration.

**Energy measurement tools.** For hardware energy measurement, we employed two widely-used tools: powercap and pynvml. Powercap [70] interfaces with RAPL (Running Average Power Limit), which is supported by both our Intel and AMD processors, to collect CPU energy consumption data through records in system files. For GPU energy monitoring, we utilized pynvml [71], the Python binding for NVIDIA Management Library (NVML), which provides direct access to GPU power sensors. Both tools offer millisecond-level sampling precision with minimal overhead to the system being measured.

**LLM selection.** In the resource-constrained environment, we select the Llama-3.1-8B-Instruct model with 8-bit quantization. In the resource-abundant environment, we deploy the Llama-3.1-8B-Instruct on one A100 and Llama-3.1-70B-Instruct on four A100.

### 4.2 Task 1: Memory Formation and Retrieval

In this experiment, we selected 22,865 tokens of ground-truth documents from LoCoMo (Internal Mem), three different sized subsets (22,226, 99,946 and 202,676 tokens) from MuSiQue (External Mem). We use LlamaIndex [72] for different indexing implementations and the retrieval top-K is set to 5. Ten experiments were repeated to take the average value.

**a. Different dataset, same number of tokens**  Table 2 shows the energy and time cost comparison of different indexing methods for different datasets with similar number of total memory token counts.

*Energy Cost Per Token.* The pattern of energy cost per token reveals two critical influencing factors. The indexing method shows substantial impact, with Vector-based (**non-LLM**) indexing demonstrating significantly lower energy

---

[2]`https://build.nvidia.com/models`

consumption ($10^{-2}$ J/token) compared to **LLM-dependent methods** (Keyword, Summary, and Graph), which require approximately $10^{+1}$ J/token. Moreover, despite similar total token counts, the memory configuration significantly affects energy efficiency, with the internal memory setup (419 smaller nodes) consistently consuming more energy than the external configuration (220 larger nodes) across all indexing methods.

*RERR and EAR.* The RERR and EAR cost metrics derive from the energy per token and essentially reflect the large differences between the different indexing approaches, with vector indexing achieving **greater per-unit retrieval effectiveness** at significantly lower energy cost ($10^{-4}$J) compared to LLM-dependent approaches ($10^{-1}$J). A distinctive pattern emerges in the cost distribution: internal data exhibits higher RERR and lower EAR values, while external memory shows the opposite trend, indicating varying energy investments required for achieving optimal Relevance and Accuracy across different scenarios. Notably, the **Graph** index, despite being the most sophisticated method, demands higher energy costs in simpler query scenarios (internal, 0.717 J and 1.035 J) compared to complex multi-hop queries (external, 0.565 J and 0.390 J), emphasizing the importance of context-appropriate index selection rather than defaulting to the most advanced option to avoid unnecessary energy expenditure.

*Time Cost.* The time cost reveals that node configuration significantly impacts both writing and reading performance, with the internal memory's more nodes requiring longer processing times for both indexing and querying compared to external data. The Vector-based method, operating without LLM, demonstrates extremely efficient processing times ($10^{-5}$ s/token) for the moderate token count (around 22k), while LLM-dependent methods demand substantially longer indexing times, highlighting the need to consider these significant preprocessing costs in specific use scenarios. Additionally, certain indexing methods show reading efficiency advantages in specific contexts, such as **Keyword**'s faster reading times (2.01e-06 s/token), while advanced methods like **Graph** indexing require much longer (4.55e-02 s/token), suggesting that optimal time efficiency requires careful matching between indexing strategies and application requirements rather than assuming more sophisticated methods will always yield better performance.

Table 2: Cost Comparison of Memory Formation (different dataset, similar number of tokens)

| Index | Internal Mem, Token Num: 22,865, Node num: **419** | | | | | External Mem, Token Num: **22,226**, Node Num: **220** | | | | |
| | RERR | EAR | Writing Time (s) | Reading Time (s) | Energy per token (J) | RERR | EAR | Writing Time (s) | Reading Time (s) | Energy per token (J) |
|---|---|---|---|---|---|---|---|---|---|---|
| Vector | **1.79e-04** | **2.12e-04** | **8.75e-05** | **3.29e-05** | 1.18e-02 | **2.25e-04** | **1.01e-04** | **5.14e-05** | 1.59e-05 | 6.12e-03 |
| Keyword | **8.15e-01** | 10.23e-01 | 3.02e-01 | 8.88e-05 | 3.27e+01 | **9.60e-01** | 5.15e-01 | 1.75e-01 | **2.01e-06** | 2.52e+01 |
| Summary | 6.21e-01 | **12.67e-01** | 3.10e-01 | 4.14e-04 | 3.40e+01 | 9.19e-01 | **5.24e-01** | 1.71e-01 | 8.77e-05 | 2.50e+01 |
| Graph | <u>7.17e-01</u> | <u>10.35e-01</u> | 3.29e-01 | **4.86e-02** | 3.63e+01 | <u>5.65e-01</u> | <u>3.90e-01</u> | 1.77e-01 | **4.55e-02** | 2.57e+01 |

**b. Same dataset, different number of tokens**   Table 3 and External Mem part in Table 2 shows the cost comparison of different indexing methods for same datasets with increasing number of memory tokens. The cost pattern analysis across increasing token numbers (from 22,226 to 99,946 to 202,676) in external memory reveals a trend in efficiency scaling. The Vector-based method shows **consistently increasing costs** across all metrics: energy per token rises from 6.12e-03 J to 9.80e-03 J, while its RERR and EAR values also increase proportionally. In contrast, LLM-dependent methods demonstrate relatively **stable cost patterns**, remaining nearly constant metrics across different token counts. The narrowing gap and the stability of the LLM-dependent approach's performance with increasing token volumes suggests its potential to remain cost-effective at scale. Nonetheless, the huge magnitude difference ($10^{-3}$ vs $10^{+1}$) in energy costs between the two still needs to be addressed first.

Table 3: Cost Comparison of Memory Formation (same dataset, different number of tokens)

| Index | External Mem, Token Number = **99,946** | | | | | External Mem, Token Num = **202,676** | | | | |
| | RERR | EAR | Writing Time (s) | Reading Time (s) | Energy per token (J) | RERR | EAR | Writing Time (s) | Reading Time (s) | Energy per token (J) |
|---|---|---|---|---|---|---|---|---|---|---|
| Vector | 2.48e-04 | 1.12e-04 | 5.16e-05 | 4.38e-05 | **8.94**e-03 | 2.53e-04 | 1.28e-04 | 5.40e-05 | 7.03e-05 | **9.80**e-03 |
| Keyword | 9.79e-01 | 6.40e-01 | 1.79e-01 | 4.76e-06 | **2.50**e+01 | 9.23e-01 | 6.59e-01 | 1.71e-01 | 9.56e-06 | **2.48**e+01 |
| Summary | 7.87e-01 | 4.94e-01 | 1.76e-01 | 2.12e-04 | **2.52**e+01 | 7.66e-01 | 4.81e-01 | 1.69e-01 | 3.73e-04 | **2.50**e+01 |
| Graph | 6.11e-01 | 3.21e-01 | 1.79e-01 | 4.29e-02 | **2.57**e+01 | 6.17e-01 | 3.23e-01 | 1.77e-01 | 3.58e-02 | **2.59**e+01 |

### 4.3 Task 2: Memory Reading, Utilization and Generation

**a. For all generations** In this experiment, we selected 100 simpler, non-multi-hop reasoning queries from the LoCoMo, 200 queries from HotpotQA & MuSiQue, and 100 academic queries from ArXiv dataset. The indexing and retrieval are fixed to the vector-based method and the retrieval top-K is set to 5. Ten experiments were repeated to take the average value. The results are shown in Table 4 and Table 5, where Table 5 is a multiples analysis based on the GERR, EFR, and ECR values in Table 4. For the choice of *query optimisation* methods, we followed the experimental results of other previous RAG evaluation work [73], choosing the HyDE method which has relatively better metrics. For the other operations, we use prompt engineering for ordinary LLM inference to implement them.

Table 4: Cost Comparison of Memory Reading, Utilization Operations

| Operation | Internal Mem. | | | | External Mem. A | | | | External Mem. B | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | GERR | EFR | ECR | latency(s) | GERR | EFR | ECR | latency(s) | GERR | EFR | ECR | latency(s) |
| Retrieval + Generation | 0.144 | 0.096 | 0.328 | 0.059 | **0.951** | 0.102 | **1.440** | 0.059 | 0.230 | **0.176** | 0.669 | **0.102** |
| + retrieval detection | 0.188 | 0.137 | 0.413 | 0.077 | **1.114** | 0.120 | **1.686** | 0.069 | 0.385 | **0.294** | 1.120 | **0.111** |
| + query optimization | 5.005 | **3.312** | 11.918 | **2.115** | **15.957** | 1.823 | **19.490** | 1.117 | 2.806 | 2.012 | 8.135 | 1.073 |
| + reranking | 0.584 | **0.413** | 1.307 | **0.245** | **3.172** | 0.356 | **4.490** | 0.202 | 0.499 | 0.361 | 1.368 | 0.192 |
| + compression | 2.542 | 1.686 | 5.261 | 1.084 | **12.884** | 2.039 | **23.653** | **1.187** | 1.193 | 0.942 | 3.403 | 0.468 |
| + all | 6.561 | **4.604** | 14.325 | **2.841** | 28.568 | 3.596 | **43.460** | 2.106 | 3.433 | 2.809 | 9.637 | 1.330 |

Table 5: Energy Cost Multiples Analysis (*Retrieval + Generation* as Baseline)

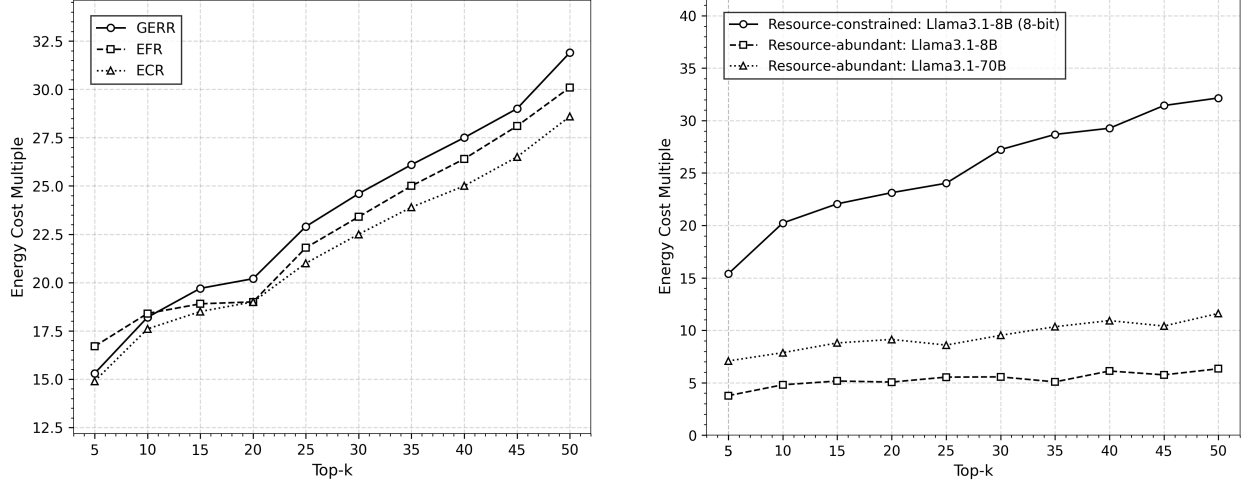| Operation | Internal Mem. | External Mem. A | External Mem. B |
|---|---|---|---|
| + retrieval detection | 1.30× | 1.17× | 1.67× |
| + query optimization | 35.19× | 15.95× | 11.93× |
| + reranking | 4.11× | 3.31× | 2.09× |
| + compression | 17.07× | 16.45× | 5.21× |
| + all | **45.70×** | **31.74×** | **15.08×** |

Note: Values represent geometric mean multiples of relative increases in GERR, EFR, and ECR

*External Mem B (ArXiv).* The arXiv dataset, characterized by longer queries and a large number of nodes (51,270), impacted retrieval and detection costs, leading to **relatively high latency** (0.102, 0.111). Longer responses contained more information and divergence points, while smaller node sizes reduced the information in individual text fragments. This mismatch made it harder for the LLM to align responses with retrieved text, leading to **higher EFR cost** (0.176, 0.294).

*External Mem A (HotpotQA & MuSiQue).* The multi-hop query introduced higher query complexity, which notably led to **high GERR and ECR metrics**. These results highlighted the additional energy costs required to improve response-query relevance and ensure response accuracy. Although the *compression* effectively alleviated the pressure of larger nodes on LLM input length, it incurred a significant energy cost (a **16-fold increase**), prompting a need to carefully weigh the cost-effectiveness of such an expensive optimization step in various scenarios.

*Internal Mem (LoCoMo).* LoCoMo consisted of simple, everyday conversational queries. Its low query complexity was reflected in the relatively **low initial GERR** (0.144) and **ECR** (0.328). However, during *query optimization*, the generated hypothetical documents for these simple tasks often introduced redundant or distracting information—particularly in time-, location-, or person-based dialogues. This "over-optimization" led to a dramatic increase in energy consumption and latency (about **36-fold rise**), underscoring the critical importance of carefully selecting operational strategies.

*Costs of different operations.* With relatively low additional costs (1.17 to 1.67-fold increase), *retrieval detection* can be acceptable as a better trade-off operation, if in scenarios with simpler queries or where the LLM's own knowledge is sufficient, timely detection can prevent unnecessary downstream operations. In contrast, sub-optimal operations can have extreme energy impacts, an example of which is *query optimisation* in internal memory. All operations (+all) resulted in energy costs increasing by **an order of magnitude or more** (45.7-fold, 31.74-fold, 15.08-fold). Although the relative increase of External Mem.B was smaller, the overhead **linearly scaled** as the memory reading workload (top-K retrieved documents) increased, as shown in Figure 4a.

(a) Cost Multiple Comparison (External Mem. B)    (b) Cost Multiple Comparison (Different Models)

Figure 4: The Relationship between Energy Cost Multiple and Memory Reading Workload (top-K value)

**b. For correct generations**    In this experiment, we selected 100 correctly answered queries each from the LoCoMo, HotpotQA and MuSiQue respectively. Through 10 experimental iterations with 100 different (input length, output length) pairs in our **resource-constrained environment**, we measured the actual energy consumption using monitoring tools and then applying OLS regression to fit the energy model 3.1. We derived the coefficients $\alpha_{K,0} = \mathbf{0.042933}$, $\alpha_{K,1} = \mathbf{9.109322}$, and $\alpha_{K,2} = \mathbf{0.000513}$ for the equation (3), where K represents Llama-3.1-8B-Instruct model with 8-bit quantization. The high R-squared value of 0.989 validates that this model effectively captures the energy-token numbers relationship.

Based on our defined minimal essential operations and data in 3.3.2, we calculated $E_g$ by substituting into equation (3) and measured the retrieval cost $E_{retr}$ and real energy cost $E_{real}$ when all operations were performed. Following Equations (6) and (7), we calculated the Generation Energy Optimality Ratio (GEOR) for the three datasets. The results show that the GEOR values are extremely low as shown in Table 6, all only **about 1-2%**. Several factors lead to excessive extra energy consumption, yielding these exceptionally low values: the retrieval and processing of many irrelevant memories, the accumulated overhead from additional operations, and the extended processing time from retrieval and computation latency. This indicate that while the LLM agent is capable of correctly answering memory-based tasks, there is a considerable gap between the system's current energy performance and the theoretical minimum, underscoring the inefficiency in energy utilization during the generation pipeline, as well as the significant room for improvement in optimising the energy efficiency of the LLM agent's design pattern.

Table 6: Generation Energy Optimality Ratio (when all operations are performed)

| Dataset | LoCoMo | HotpotQA | MuSiQue |
|---------|--------|----------|---------|
| GEOR | 1.378% | 1.012% | 1.456% |

## 4.4    Task 3: Resource-constrained vs. Resource-abundant Environment

In this part we repeated the same experiments for Task 1 and 2 in resource-abundant environment using Llama-3.1-8B-Instruct and Llama-3.1-70B-Instruct. The results shown in Figure 5 represent the average metrics of the three LLM-based index methods (Keyword, Summary, and Graph) for the Memory Formation and Retrieval phases, and Figure 4b, the energy of the Generation phases for performing all operations on the Arxiv dataset compared to the multiples of energy consumption of the baseline (Retrieval + Generation).

*Cost-Effectiveness Paradox.* Despite operating with reduced precision (8-bit quantization), the resource-constrained environment shows notably poor cost-efficiency ratios. The resource-abundant environment, running full-precision 8B model, achieves lower RERR (0.642 vs 0.769) and EAR (0.352 vs 0.488), and consumed less energy per memory token (18.7 vs 25.2 joules). Given that RERR and EAR measure the energy cost required to achieve unit retrieval relevance and accuracy respectively, this counter-intuitive result reveals that resource-constrained environments are forced to pay higher operational costs (35% more energy per token) while achieving lower effectiveness.
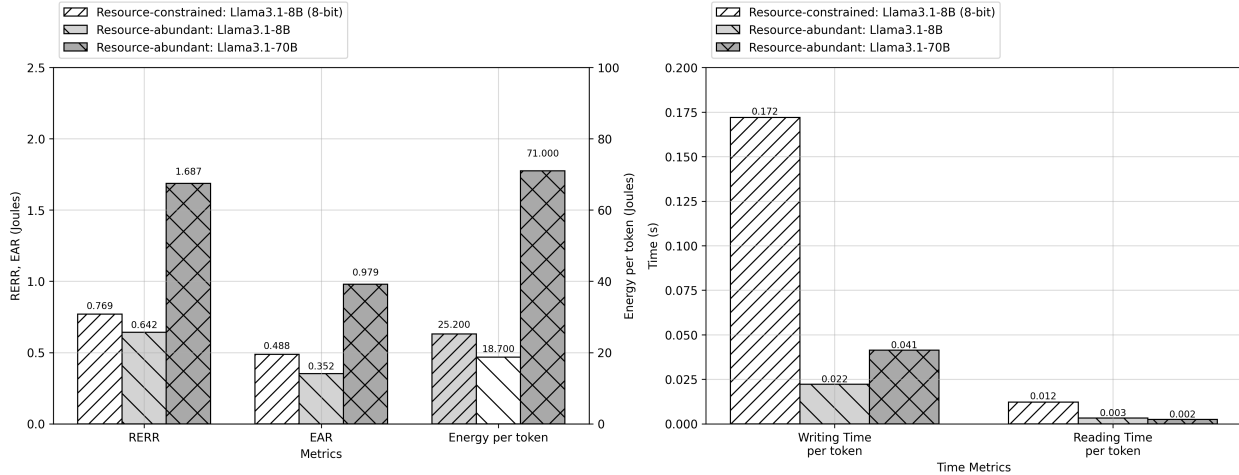
13

Figure 5: Comparison of Memory Formation in Resource-constrained vs. Resource-abundant Environments

*Model Size Scaling Benefits.* The 70B model consumes approximately three times more energy per memory token than the quantized 8B model (67.5 vs 25.2 joules), while its RERR (1.687 vs 0.769) and EAR (0.979 vs 0.488) are only about twice as high. This disproportionate relationship indicates that the 70B model achieves higher relevance and accuracy improvements relative to its energy consumption increase. This advantageous scaling property - where performance improvements outpace energy cost increases - is exclusively available to resource-abundant environments.

*Service Quality Disparity.* The most striking disparity appears in processing efficiency. The resource-constrained environment exhibits dramatically slower indexing speeds, with writing time per token (0.172s) being approximately 8 times slower than the resource-abundant environment (0.022s for 8B model), while the difference in read time is also four to six times greater. This difference in processing speed signals that resource-constrained environments are paying for longer preprocessing times but gaining a more sub-optimal quality of service, potentially rendering these LLM agent systems impractical for real-world applications in resource-limited settings.

*Optimization Penalty Paradox.* Further analysis of the Generation phase (Figure 4b) reveals an even more severe energy efficiency gap with optimization operations (reranking, compression, etc.). While all configurations show increasing energy costs as the workload grows, the resource-constrained environment exhibits a dramatically steeper curve - from 15x to 32x baseline consumption as top-K increases from 5 to 50. In comparison, resource-abundant environments show more moderate increases (from 4x to 6x for 8B, and 7x to 12x for 70B). This divergent scaling behavior indicates that although the baseline energy overhead might be lower in resource-constrained environments, their efficiency deteriorates much more rapidly with increased context length and additional processing steps. This creates a paradoxical situation where attempts to improve system performance through optimization operations actually widen the energy efficiency gap between resource-abundant and resource-constrained deployments.

## 5 Discussions and Conclusions

Our experimental investigation into LLM agent systems provides empirical evidence for the Sustainable AI Trilemma and its constituent dilemmas, revealing how the tensions between capability, digital equity, and environmental sustainability manifest in practical deployments. Through a systematic examination of memory operations in LLM agents, we observe several critical patterns that validate and deepen our understanding of these challenges.

The **Operational Dilemma** manifests through stark energy consumption patterns in LLM agent systems. The fundamental tension between **scalability and sustainability** is evident in the orders-of-magnitude difference between LLM-dependent and non-LLM methods, with optimization attempts often leading to significant energy penalties. The remarkably low Generation Energy Optimality Ratio (GEOR) of 1-2% highlights how current architectures operate far from optimal efficiency, suggesting that without fundamental breakthroughs in energy efficiency, the widespread deployment of LLM agents may face significant sustainability barriers.

The **Evaluation Dilemma** reflects the crucial transition from **mere awareness to precise measurement** of AI's environmental impact. While previous approaches relied on simple proxies like the product of GPU power and runtime, our evaluation framework introduces targeted metrics and more concrete evidences that reveal nuanced trade-offs

14

and previously obscured patterns between energy costs and system performance. This progress in measurement capability, though significant, also highlights the continuing challenge of developing comprehensive assessment tools for increasingly complex AI applications.

The **Design Dilemma** challenges the prevailing "LLM-centric" paradigm that delegates complete operational control to LLM. Our findings reveal the inefficiencies and overhead associated with **fully LLM autonomous** operation while only some LLM operations can be **cost-effective** like retrieval detection. The systematic increase in energy costs with task complexity suggests the need for a more balanced approach that strategically combines LLM and traditional low-cost non-LLM methods. In addition, whether LLMs can accurately assess the boundaries of their capabilities [39], adaptively avoid inappropriate operations and perform optimal operation allocation are all key issues that need to be further explored.

The **Access Dilemma** emerges as a profound manifestation of digital inequality, where the relationship between resources **barriers and benefits** creates a negative feedback loop of technological disparity: limited computational resources lead to both higher operational costs and reduced capabilities but also experience disproportionate penalties when attempting to optimize their systems. This troubling "technological poverty trap" creates a fundamental barrier to AI democratization, where the attempts to improve accessibility through techniques like model quantization and RAG optimization may widen the capability gap between resource-abundant and resource-constrained deployments, effectively preventing disadvantaged communities from realizing the transformative potential of AI technologies. This pattern suggests that the current trajectory of AI development may inadvertently amplify existing digital divide rather than bridge them.

This work has highlighted the complex, interconnected challenges in developing sustainable AI systems, particularly for LLM agents. Our findings demonstrate that addressing any single aspect of sustainability often creates ripple effects across technical, social, and environmental dimensions, necessitating a holistic approach to solution design. Looking ahead, our research will focus on developing practical frameworks for sustainable AI practices that can effectively reduce environmental impact while preserving the utility of LLM agents. Through continued investigation of responsible AI design principles and their implementation, we aim to enable the equitable deployment of these systems across diverse contexts, including resource-constrained environments. This future direction emphasizes the critical balance between advancing AI capabilities and ensuring their long-term sustainability through coordinated technical innovation and thoughtful policy development.

## Acknowledgments

## References

[1] David Patterson et al. "The Carbon Footprint of Machine Learning Training Will Plateau, Then Shrink". In: *Computer* 55.7 (2022), pp. 18–28. DOI: 10.1109/MC.2022.3148714.

[2] Andrew A Chien et al. "Reducing the Carbon Impact of Generative AI Inference (today and in 2035)". In: *Proceedings of the 2nd Workshop on Sustainable Computer Systems*. HotCarbon '23. Boston, MA, USA: Association for Computing Machinery, 2023. DOI: 10.1145/3604930.3605705.

[3] Radosvet Desislavov, Fernando Martínez-Plumed, and José Hernández-Orallo. "Trends in AI inference energy consumption: Beyond the performance-vs-parameter laws of deep learning". In: *Sustainable Computing: Informatics and Systems* 38 (2023), p. 100857. ISSN: 2210-5379. DOI: https://doi.org/10.1016/j.suscom.2023.100857.

[4] Peng Jiang, Christian Sonne, Wangliang Li, Fengqi You, and Siming You. "Preventing the Immense Increase in the Life-Cycle Energy and Carbon Footprints of LLM-Powered Intelligent Chatbots". In: *Engineering* 40 (2024), pp. 202–210. ISSN: 2095-8099. DOI: https://doi.org/10.1016/j.eng.2024.04.002.

[5] Tong Wu et al. "Triple Bottom Line or Trilemma? Global Tradeoffs Between Prosperity, Inequality, and the Environment". In: *World Development* 178 (2024), p. 106595.

[6] Zhiheng Xi et al. "The Rise and Potential of Large Language Model Based Agents: A Survey". 2023. arXiv: 2309.07864 [cs.AI].

[7] Sunder Ali Khowaja, Parus Khuwaja, Kapal Dev, Weizheng Wang, and Lewis Nkenyereye. "Chatgpt needs spade (sustainability, privacy, digital divide, and ethics) evaluation: A review". In: *Cognitive Computation* (2024), pp. 1–23. DOI: 10.1007/s12559-024-10285-1.

[8]    Ahmad Faiz et al. "LLMCarbon: Modeling the end-to-end Carbon Footprint of Large Language Models". 2024. arXiv: 2309.14393 [cs.CL].

[9]    Ekkehard Ernst. "The AI trilemma: Saving the planet without ruining our jobs". In: *Frontiers in Artificial Intelligence* 5 (2022). ISSN: 2624-8212. DOI: 10.3389/frai.2022.886561.

[10]   Sophia Falk and Aimee van Wynsberghe. "Challenging AI for Sustainability: what ought it mean?" In: *AI and Ethics* 4.4 (2024), pp. 1345–1355.

[11]   Zhenghua Chen, Min Wu, Alvin Chan, Xiaoli Li, and Yew-Soon Ong. "Survey on AI Sustainability: Emerging Trends on Learning Algorithms and Research Challenges". In: *IEEE Computational Intelligence Magazine* 18.2 (2023), pp. 60–77. DOI: 10.1109/MCI.2023.3245733.

[12]   Alexandre Lacoste, Alexandra Luccioni, Victor Schmidt, and Thomas Dandres. "Quantifying the Carbon Emissions of Machine Learning". 2019. arXiv: 1910.09700 [cs.CY].

[13]   Loïc Lannelongue, Jason Grealey, and Michael Inouye. "Green Algorithms: Quantifying the Carbon Footprint of Computation". In: *Advanced Science* 8.12 (2021), p. 2100707. DOI: https://doi.org/10.1002/advs.202100707.

[14]   Yue Liu et al. "Agent design pattern catalogue: A collection of architectural patterns for foundation model based agents". In: *Journal of Systems and Software* 220 (2025), p. 112278. ISSN: 0164-1212. DOI: 10.1016/j.jss.2024.112278.

[15]   Himel Ghosh. "Enabling Efficient Serverless Inference Serving for LLM (Large Language Model) in the Cloud". 2024. arXiv: 2411.15664 [cs.DC].

[16]   Josu Diaz-de-Arcaya, Ana I. Torre-Bastida, Gorka Zárate, Raúl Miñón, and Aitor Almeida. "A Joint Study of the Challenges, Opportunities, and Roadmap of MLOps and AIOps: A Systematic Survey". In: *ACM Comput. Surv.* 56.4 (Oct. 2023). ISSN: 0360-0300. DOI: 10.1145/3625289.

[17]   Renzhe Yu, Zhen Xu, Sky CH-Wang, and Richard Arum. "Whose ChatGPT? Unveiling Real-World Educational Inequalities Introduced by Large Language Models". 2024. arXiv: 2410.22282 [cs.CY].

[18]   Rehema Baguma, Hajarah Namuwaya, Joyce Nakatumba-Nabende, and Qazi Mamunur Rashid. "Examining Potential Harms of Large Language Models (LLMs) in Africa". In: *International Conference on Safe, Secure, Ethical, Responsible Technologies and Emerging Applications*. Springer. 2023, pp. 3–19.

[19]   Yuheng Cheng et al. "Exploring Large Language Model based Intelligent Agents: Definitions, Methods, and Prospects". 2024. arXiv: 2401.03428 [cs.AI].

[20]   Xu Huang et al. "Understanding the planning of LLM agents: A survey". 2024. arXiv: 2402.02716 [cs.AI].

[21]   Zhuocheng Shen. "LLM With Tools: A Survey". 2024. arXiv: 2409.18807 [cs.AI].

[22]   Siyun Zhao et al. "Retrieval Augmented Generation (RAG) and Beyond: A Comprehensive Survey on How to Make your LLMs use External Data More Wisely". 2024. arXiv: 2409.14924 [cs.CL].

[23]   Zeyu Zhang et al. "A Survey on the Memory Mechanism of Large Language Model based Agents". 2024. arXiv: 2404.13501 [cs.AI].

[24]   Yue Zhang et al. "Siren's Song in the AI Ocean: A Survey on Hallucination in Large Language Models". 2023. arXiv: 2309.01219 [cs.CL].

[25]   Yi Ding and Tianyao Shi. "Sustainable LLM Serving: Environmental Implications, Challenges, and Opportunities : Invited Paper". In: *2024 IEEE 15th International Green and Sustainable Computing Conference (IGSC)*. 2024, pp. 37–38. DOI: 10.1109/IGSC64514.2024.00016.

[26]   Siddharth Samsi et al. "From Words to Watts: Benchmarking the Energy Costs of Large Language Model Inference". In: *2023 IEEE High Performance Extreme Computing Conference (HPEC)*. 2023, pp. 1–9. DOI: 10.1109/HPEC58863.2023.10363447.

[27]   Grant Wilkins, Srinivasan Keshav, and Richard Mortier. "Offline Energy-Optimal LLM Serving: Workload-Based Energy Models for LLM Inference on Heterogeneous Systems". 2024. arXiv: 2407.04014 [cs.DC].

[28]   Mauricio Fadel Argerich and Marta Patiño-Martínez. "Measuring and Improving the Energy Efficiency of Large Language Models Inference". In: *IEEE Access* 12 (2024), pp. 80194–80207. DOI: 10.1109/ACCESS.2024.3409745.

[29]   Tina Vartziotis et al. "Learn to Code Sustainably: An Empirical Study on LLM-based Green Code Generation". 2024. arXiv: 2403.03344 [cs.SE].

[30]   Paul Delanoë, Dieudonné Tchuente, and Guillaume Colin. "Method and evaluations of the effective gain of artificial intelligence models for reducing CO2 emissions". In: *Journal of Environmental Management* 331 (2023), p. 117261. ISSN: 0301-4797. DOI: https://doi.org/10.1016/j.jenvman.2023.117261.

[31]   Zhi Jing, Yongye Su, and Yikun Han. "When Large Language Models Meet Vector Databases: A Survey". 2024. arXiv: 2402.01763 [cs.DB].

[32] Junru Lu et al. "MemoChat: Tuning LLMs to Use Memos for Consistent Long-Range Open-Domain Conversation". 2023. arXiv: 2308.08239 [cs.CL].

[33] Wanjun Zhong, Lianghong Guo, Qiqi Gao, He Ye, and Yanlin Wang. "Memorybank: Enhancing large language models with long-term memory". In: *Proceedings of the AAAI Conference on Artificial Intelligence*. Vol. 38. 2024, pp. 19724–19731.

[34] Bing Wang et al. "Enhancing Large Language Model with Self-Controlled Memory Framework". 2024. arXiv: 2304.13343 [cs.CL].

[35] Charles Packer et al. "MemGPT: Towards LLMs as Operating Systems". 2024. arXiv: 2310.08560 [cs.AI].

[36] Gibbeum Lee, Volker Hartmann, Jongho Park, Dimitris Papailiopoulos, and Kangwook Lee. "Prompted LLMs as Chatbot Modules for Long Open-domain Conversation". In: *Findings of the Association for Computational Linguistics: ACL 2023*. Ed. by Anna Rogers, Jordan Boyd-Graber, and Naoaki Okazaki. Toronto, Canada: Association for Computational Linguistics, July 2023, pp. 4536–4554. DOI: 10.18653/v1/2023.findings-acl.277.

[37] Lei Liu et al. "Think-in-Memory: Recalling and Post-thinking Enable LLMs with Long-Term Memory". 2023. arXiv: 2311.08719 [cs.CL].

[38] Ali Modarressi, Ayyoob Imani, Mohsen Fayyaz, and Hinrich Schütze. "RET-LLM: Towards a General Read-Write Memory for Large Language Models". 2024. arXiv: 2305.14322 [cs.CL].

[39] Ruiyang Ren et al. "Investigating the Factual Knowledge Boundary of Large Language Models with Retrieval Augmentation". 2024. arXiv: 2307.11019 [cs.CL].

[40] Hanxing Ding, Liang Pang, Zihao Wei, Huawei Shen, and Xueqi Cheng. "Retrieve Only When It Needs: Adaptive Retrieval Augmentation for Hallucination Mitigation in Large Language Models". 2024. arXiv: 2402.10612 [cs.CL].

[41] Yile Wang, Peng Li, Maosong Sun, and Yang Liu. "Self-Knowledge Guided Retrieval Augmentation for Large Language Models". In: *Findings of the Association for Computational Linguistics: EMNLP 2023*. Ed. by Houda Bouamor, Juan Pino, and Kalika Bali. Singapore: Association for Computational Linguistics, Dec. 2023, pp. 10303–10315. DOI: 10.18653/v1/2023.findings-emnlp.691.

[42] Denny Zhou et al. "Least-to-Most Prompting Enables Complex Reasoning in Large Language Models". 2023. arXiv: 2205.10625 [cs.AI].

[43] Jason Wei et al. "Chain-of-Thought Prompting Elicits Reasoning in Large Language Models". 2023. arXiv: 2201.11903 [cs.CL].

[44] Ofir Press et al. "Measuring and Narrowing the Compositionality Gap in Language Models". In: *Findings of the Association for Computational Linguistics: EMNLP 2023*. Singapore: Association for Computational Linguistics, Dec. 2023, pp. 5687–5711. DOI: 10.18653/v1/2023.findings-emnlp.378.

[45] LangChain. "LangChain MultiQueryRetriever Documentation". Accessed: 2025-01-10. 2024. URL: https://python.langchain.com/v0.1/docs/modules/data_connection/retrievers/MultiQueryRetriever/.

[46] Luyu Gao, Xueguang Ma, Jimmy Lin, and Jamie Callan. "Precise Zero-Shot Dense Retrieval without Relevance Labels". In: *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. Ed. by Anna Rogers, Jordan Boyd-Graber, and Naoaki Okazaki. Toronto, Canada: Association for Computational Linguistics, July 2023, pp. 1762–1777. DOI: 10.18653/v1/2023.acl-long.99.

[47] Huaixiu Steven Zheng et al. "Take a Step Back: Evoking Reasoning via Abstraction in Large Language Models". 2024. arXiv: 2310.06117 [cs.LG].

[48] Xinbei Ma, Yeyun Gong, Pengcheng He, Hai Zhao, and Nan Duan. "Query Rewriting in Retrieval-Augmented Large Language Models". In: *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*. Singapore: Association for Computational Linguistics, Dec. 2023, pp. 5303–5315. DOI: 10.18653/v1/2023.emnlp-main.322.

[49] Stephen Robertson and Hugo Zaragoza. "The probabilistic relevance framework: BM25 and beyond". In: *Foundations and Trends in Information Retrieval* 3.4 (2009), pp. 333–389.

[50] Nils Reimers and Iryna Gurevych. "Sentence-BERT: Sentence Embeddings using Siamese BERT-Networks". 2019. arXiv: 1908.10084 [cs.CL].

[51] Shengyao Zhuang, Bing Liu, Bevan Koopman, and Guido Zuccon. "Open-source Large Language Models are Strong Zero-shot Query Likelihood Models for Document Ranking". 2023. arXiv: 2310.13243 [cs.IR].

[52] Xueguang Ma, Xinyu Zhang, Ronak Pradeep, and Jimmy Lin. "Zero-Shot Listwise Document Reranking with a Large Language Model". 2023. arXiv: 2305.02156 [cs.IR].

[53] Devendra Sachan et al. "Improving Passage Retrieval with Zero-Shot Question Generation". In: *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*. Abu Dhabi, United Arab Emirates: Association for Computational Linguistics, Dec. 2022, pp. 3781–3797. DOI: `10.18653/v1/2022.emnlp-main.249`.

[54] Weiwei Sun et al. "Is ChatGPT Good at Search? Investigating Large Language Models as Re-Ranking Agents". In: *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*. Ed. by Houda Bouamor, Juan Pino, and Kalika Bali. Singapore: Association for Computational Linguistics, Dec. 2023, pp. 14918–14937. DOI: `10.18653/v1/2023.emnlp-main.923`.

[55] Zhen Qin et al. "Large Language Models are Effective Text Rankers with Pairwise Ranking Prompting". In: *Findings of the Association for Computational Linguistics: NAACL 2024*. Ed. by Kevin Duh, Helena Gomez, and Steven Bethard. Mexico City, Mexico: Association for Computational Linguistics, June 2024, pp. 1504–1518. DOI: `10.18653/v1/2024.findings-naacl.97`.

[56] Percy Liang et al. "Holistic Evaluation of Language Models". 2023. arXiv: 2211.09110 [`cs.CL`].

[57] Haoyan Yang et al. "PRCA: Fitting Black-Box Large Language Models for Retrieval Question Answering via Pluggable Reward-Driven Contextual Adapter". 2023. arXiv: 2310.18347 [`cs.CL`].

[58] Fangyuan Xu, Weijia Shi, and Eunsol Choi. "RECOMP: Improving Retrieval-Augmented LMs with Compression and Selective Augmentation". 2023. arXiv: 2310.04408 [`cs.CL`].

[59] Yucheng Li, Bo Dong, Chenghua Lin, and Frank Guerin. "Compressing Context to Enhance Inference Efficiency of Large Language Models". 2023. arXiv: 2310.06201 [`cs.CL`].

[60] Huiqiang Jiang, Qianhui Wu, Chin-Yew Lin, Yuqing Yang, and Lili Qiu. "LLMLingua: Compressing Prompts for Accelerated Inference of Large Language Models". 2023. arXiv: 2310.05736 [`cs.CL`].

[61] Zhuoshi Pan et al. "LLMLingua-2: Data Distillation for Efficient and Faithful Task-Agnostic Prompt Compression". 2024. arXiv: 2403.12968 [`cs.CL`].

[62] Hao Yu et al. "Evaluation of Retrieval-Augmented Generation: A Survey". 2024. arXiv: 2405.07437 [`cs.CL`].

[63] Adyasha Maharana et al. "Evaluating Very Long-Term Conversational Memory of LLM Agents". 2024. arXiv: 2402.17753 [`cs.CL`].

[64] Zhilin Yang et al. "HotpotQA: A Dataset for Diverse, Explainable Multi-hop Question Answering". 2018. arXiv: 1809.09600 [`cs.CL`].

[65] Harsh Trivedi, Niranjan Balasubramanian, Tushar Khot, and Ashish Sabharwal. "MuSiQue: Multihop Questions via Single-hop Question Composition". In: *Transactions of the Association for Computational Linguistics* 10 (2022), pp. 539–554.

[66] James Calam. "Ai arxiv dataset". Accessed: 2025-01-10. 2023. URL: `https://huggingface.co/datasets/jamescalam/ai-arxiv`.

[67] Matouš Eibich, Shivay Nagpal, and Alexander Fred-Ojala. "ARAGOG: Advanced RAG Output Grading". 2024. arXiv: 2404.01037 [`cs.CL`].

[68] UpTrain AI. "UpTrain". Accessed: 2025-01-10. 2024. URL: `https://github.com/uptrain-ai/uptrain`.

[69] Chin-Yew Lin. "ROUGE: A Package for Automatic Evaluation of Summaries". In: *Text Summarization Branches Out*. Barcelona, Spain: Association for Computational Linguistics, July 2004, pp. 74–81. URL: `https://aclanthology.org/W04-1013/`.

[70] "Powercap Linux Kernel Interface". Accessed: 2025-01-10. 2024. URL: `https://www.kernel.org/doc/html/latest/power/powercap/powercap.html`.

[71] NVIDIA Corporation. "NVIDIA Management Library (NVML) Python Bindings". Accessed: 2025-01-10. 2024. URL: `https://pypi.org/project/nvidia-ml-py/`.

[72] Jerry Liu. "LlamaIndex". Nov. 2022. URL: `https://github.com/jerryjliu/llama_index`.

[73] Xiaohua Wang et al. "Searching for best practices in retrieval-augmented generation". In: *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*. 2024, pp. 17716–17736.