



A Learnable Agent Collaboration Network Framework for Personalized Multimodal AI Search Engine

Yunxiao Shi
Yunxiao.Shi@student.uts.edu.au
University of Technology Sydney
Sydney, NSW, Australia

Min Xu*
Min.Xu@uts.edu.au
University of Technology Sydney
Sydney, NSW, Australia

Haimin Zhang
Haimin.Zhang@uts.edu.au
University of Technology Sydney
Sydney, NSW, Australia

Xing Zi
Xing.Zi-1@student.uts.edu.au
University of Technology Sydney
Sydney, NSW, Australia

Qiang Wu
Qiang.Wu@uts.edu.au
University of Technology Sydney
Sydney, NSW, Australia

Abstract

Large language models (LLMs) and retrieval-augmented generation (RAG) techniques have revolutionized traditional information access, enabling AI agent to search and summarize information on behalf of users during dynamic dialogues. Despite their potential, current AI search engines exhibit considerable room for improvement in several critical areas. These areas include the support for multimodal information, the delivery of personalized responses, the capability to logically answer complex questions, and the facilitation of more flexible interactions. This paper proposes a novel AI Search Engine framework called the Agent Collaboration Network (ACN). The ACN framework consists of multiple specialized agents working collaboratively, each with distinct roles such as Account Manager, Solution Strategist, Information Manager, and Content Creator. This framework integrates mechanisms for picture content understanding, user profile tracking, and online evolution, enhancing the AI search engine's response quality, personalization, and interactivity. A highlight of the ACN is the introduction of a Reflective Forward Optimization method (RFO), which supports the online synergistic adjustment among agents. This feature endows the ACN with online learning capabilities, ensuring that the system has strong interactive flexibility and can promptly adapt to user feedback. This learning method may also serve as an optimization approach for agent-based systems, potentially influencing other domains of agent applications.

CCS Concepts

• Information systems → Web applications; • Computing methodologies → Multi-agent systems; Multi-agent planning; Information extraction.

*Corresponding author.



This work is licensed under a Creative Commons Attribution-NonCommercial International 4.0 License.

MMGR '24, October 28, 2024, Melbourne, VIC, Australia.

© 2024 Copyright held by the owner/author(s).

ACM ISBN 979-8-4007-1202-9/24/10

<https://doi.org/10.1145/3689091.3690087>

Keywords

Multimodal, Information Retrieval and Generation, Personalized Search, Multi-agent System

ACM Reference Format:

Yunxiao Shi, Min Xu, Haimin Zhang, Xing Zi, and Qiang Wu. 2024. A Learnable Agent Collaboration Network Framework for Personalized Multimodal AI Search Engine. In *Proceedings of the 2nd International Workshop on Deep Multimodal Generation and Retrieval (MMGR '24)*, October 2, 2024, Melbourne, VIC, Australia. ACM, New York, NY, USA, 9 pages. <https://doi.org/10.1145/3689091.3690087>

1 Introduction

In today's information-saturated world, information retrieval systems play a crucial role in sifting through vast data to find content that resonates with individual needs, thereby alleviating the problem of information overload. For years, traditional search engines like Google and Bing have been the primary tools for this task. However, recent advancements in large language models (LLMs) [6, 8] and retrieval-augmented generation (RAG) techniques [4, 16, 23] have given rise to a new generation of AI-powered search engines, such as Perplexity and Tiangong. These innovations have revolutionized information access by shifting from static query inputs to interactive dialogues with AI agents. Instead of manually browsing through multiple web pages, users can now rely on AI agents to synthesize and present the most relevant information to meet the information gaining requirements.

However, current AI search engines still have several aspects that need improvement. (1) *Multimodal Information Support*. Existing AI search engines primarily generate pure text content, whereas web content encompasses various modalities including text, images, tables, and videos [13]. The support of multimodal content understanding is essential for yielding high response quality and rich presentation content. (2) *Personalized Response*. Current AI search engines deliver uniform content to different users, overlooking the key factor of personalization and customization. While traditional search engines have incorporated some personalization features [18, 25, 33], AI search engines have yet to effectively integrate this aspect. For instance, when I asked GPT-4, Perplexity for muscle-building diet recommendations as an Indian, they all suggested beef as a primary protein source, which contradicts the cultural and dietary restrictions of Indians. (3) *Answering Complex Logic*

Requirement. Current AI search engines can handle simple information retrieval and generation tasks but struggle with complex, logic-intensive queries. Such queries often require multi-keyword searches and iterative retrieval processes, and the generated information needs logical coherence and strategic formulation. (4) *Timely Learning and Adjustment.* Current AI agents are "expert-centric", relying heavily on pre-set prompts and workflows [32], restricting their ability to autonomously adapt based on users' feedback.

Motivated by the aim of addressing these limitations, we propose a AI search engine framework named the Agent Collaboration Network (ACN). This framework comprises multiple agents, each performing distinct roles, including Account Manager, Solution Strategist, Information Manager, and Content Creator. The Account Manager interacts with users, tracks user profiles, gathers feedback, and transfer user information searching requirements to the Solution Strategist. The Solution Strategist takes account the user profile, uses a chain of thought method to solve complex user requirement step-by-step, and logically plans the article outline. It allocates information retrieval tasks to the Information Manager and content generation tasks to the Content Creator. Once everything is ready, the Solution Strategist triggers the Finalize Article action, completing the generation of multimodal content and delivering the results to the Account Manager. The Information Manager handles multimodal information retrieval, while the Content Creator generates multimodal content tailored to specific users based on the Solution Strategist's instruction and user profiles.

Additionally, we have designed an optimization algorithm named Reflective Forward Optimization (RFO) for the ACN, which can automatically adjust based on user feedback. We first design a LLM-based optimizer, which can inspect intermediate results within the agent-to-agent workflow and generates reflective reviews based on a given feedback. These reviews help improve adjustable parameters such as agent prompts, function parameters, and system settings while providing further feedback to the called agent. By running the RFO algorithm along the response-generating agent call stack, we obtain a collection of reviews for each agent. We then aggregate all review suggestions and use the LLM to update each agent, ultimately refining the entire ACN. This timely online-learning method enhances the flexibility of interactions and aligns the ACN more closely with the user's requirement.

In summary, our work provides several key contributions:

- We propose a novel AI search engine framework named Agent Collaboration Network (ACN), which incorporates a specially designed agent-learning method called Reflective Forward Optimization (RFO). The ACN surpasses traditional AI search engines by supporting multimodal content output, personalized content generation, and the creation of more logically structured and complex information. Additionally, it can continuously adjust and learn based on user feedback promptly.
- We design a synthetic dataset and use LLM-played judge to verify the effectiveness of ACN compared to SOTA Tian-Gong and Perplexity AI search engines, demonstrating its superior ability to generate engaging information with multimodality, logical-well, useful content, provide personalized user experience.

- We point out the current research gap in evaluating AI search engines' responsiveness to user feedback. We have analyzed the feasibility of experiments and outlined future plans to address this deficiency.

2 Related Works

2.1 AI Search Engine

AI search engines represent the convergence of large language models (LLMs), retrieval-augmented generation (RAG), and intelligent agent technologies, heralding a new era of search engine innovation. Given the nascent stage of AI search engine technology, we categorize the information retrieval process into six distinct phases [12]: identification of information retrieval requirements [14], retrieval augmentation [19], information retrieval and knowledge gathering [29], knowledge caching [30], knowledge filtering and ranking [5, 28], and LLM-based content generation, followed by verification and refinement [10]. Certain AI search engines do not adopt a conversational interface, requiring users to input search queries into a search bar, thereby omitting the initial phase of identifying information retrieval needs. Retrieval enhancement can be selectively integrated before the information retrieval and knowledge gathering phase, generating multiple search keywords to ensure a higher recall rate of relevant knowledge. Knowledge caching post-retrieval mitigates resource consumption associated with the retrieval and gathering phase, enhancing system responsiveness and efficiency. Knowledge filtering and ranking technique is employed before the generation phase, eliminate irrelevant information, thereby enhancing the precision of the retrieved knowledge, improving the robustness of LLM responses, and ensuring content quality. The verification and refinement phase post-generation ensures the factual accuracy of the AI-generated content and optimizes its presentation format for the user. Therefore, among these six phases, the steps of information retrieval and knowledge gathering, as well as LLM-based content generation, are necessary, corresponding to the retrieve-then-read pipeline of RAG. The other steps are optional according to practical application scenarios.

The advanced studies mentioned above focus primarily on ensuring the precision and efficiency of information generation in AI search engines. Beyond these critical metrics, our study more emphasizes content richness, personalization, and interactivity. These factors are crucial for maintaining the attractiveness of the content and enhancing the overall user experience with AI search engines.

2.2 Personalized Generation of LLM

Recent studies [15] emphasize the importance of personalizing large language models (LLMs) beyond aggregate fine-tuning methods like RLHF, as these may not fully capture diverse user preferences and values. Micro-level preference learning can better align models with individual users. Current personalization techniques mainly involve prompt tuning, which models user profiles based on historical search data to prompt LLM generating tailored outputs. Basic approaches use the entire user action history for prompting, while more advanced methods selectively retrieve relevant user data using memory mechanisms [34]. To address potential information loss,

[21] proposes a task-aware user profile summarization for prompting. Another approach [3] constructs knowledge graphs from user search and browsing activities to enhance prompt relevance.

While existing work focuses on prompt construction and user profiling, our Agent Collaboration Network (ACN) architecture shifts the emphasis to tuning agents across each step of the AI search engine workflow, aiming for a more integrated and efficient personalization strategy.

2.3 AI Agent

AI agents can interact with environments, dynamically selecting optimal actions to achieve predetermined objectives. With the advent of LLMs capable of function calls [20, 27], AI agents driven by such models exhibit exceptional intelligence and flexibility. This advancement paves the way for the evolution of RAG technology towards an Agentic RAG paradigm¹. Recent studies have introduced frameworks where multiple interconnected agents collaborate to support a wide array of tasks [11, 17]. Research indicates that agents can enhance their capabilities through reflective thinking [7] or by leveraging optimization techniques analogous to neural networks [32], highlighting significant potential for online, training-free self-adjustment of Agent-based applications.

Our work pioneers a universal framework for AI search engines utilizing multi-agent collaboration named ACN. We also introduce an optimization method enabling real-time learning and adaptation based on user feedback. Such adaptability enhances the personalized and interactive capabilities of AI search engines, offering a more tailored and responsive user experience.

3 Agent Collaboration Network

The proposed ACN framework comprises multiple agents, include the Account Manager, Solution Strategist, Information Manager, and Content Creator. Each one with distinct roles and responsibilities, and they collaborate dynamically to deliver satisfying response. The framework with a case study is presented in Figure 1.

3.1 Account Manager

The Account Manager Agent plays a pivotal role in engaging with users to comprehend their needs, monitor their interests, and assist them in articulating precise requirements. The Account Manager also serves as a critical communication conduit to Solution Strategist agent. When confirming the users' specific information retrieval needs, it conveys detailed user requirement to the Solution Strategist agent.

The Account Manager continuously tracks users' interest and information for building their profile. For a given user u , let $P = \{d_1, d_2, \dots, d_{|P|}\}$ represent the user profile, where each element $d_i = (\text{text}, \text{attitude}) \in P$ encapsulates a concise description of the user's fundamental information or interest preferences. For d_i detailing basic information, the attitude is labeled as 'None', whereas for interest preferences, the attitude is categorized into {Positive, Neutral, Negative}. We use d_{new} to denote the captured profile description during the Account Manager's ongoing interactions with the user. To integrate this new description into the pre-existing set P , we assess the topic similarity between d_{new}

and each d_i , represented as $S(d_{\text{new}}, d_i)$. This assessment leverages the bge-m3 model, which generates both dense and sparse embeddings—also referred to as lexical weights—for the text of each description. This dual embedding approach enables a hybrid similarity computation, effectively balancing keyword matching with semantic alignment between descriptions. We establish a similarity threshold γ ; if $S(d_{\text{new}}, d_i) \geq \gamma$, d_{new} replaces d_i , otherwise, d_{new} is appended to the set P as a new element.

The user's feedback can trigger the Account Manager Agent's function of Accepting Feedback and Reflection, the feedback will then be conveyed to the ACN optimizer for improving the collaborations of ACN. The details are described in Section 4.

3.2 Solution Strategist

The Solution Strategist Agent adhere to the detailed requirements from Account Manager Agent, and take account users' profile, then meticulously plans and orchestrates the process of addressing the information retrieval and generation task. By leveraging LLMs and employing the chain of thought method, the agent generates a structured and detailed pathway for problem-solving. As part of its strategic plan, the Solution Strategist Agent is empowered to execute specific actions such as Search Information, Generate Content, and Finalize Article.

For the action Search Information, the Solution Strategist Agent allocates a retrieval task to the Information Manager Agent, providing a precise search query. In the case of Generate Content, a text generation task is assigned to the Content Creator Agent, accompanied by a detailed creation requirement. Finally, for the action Finalize Article, the Solution Strategist Agent merges the generated content and delivers it back to the Account Manager Agent.

3.3 Information Manager

The Information Manager Agent is dedicated to retrieving pertinent information, utilizing the Bing Search Engine v7 API to access real-time and up-to-date web content. This agent retrieves and converts webpage content into markdown format, ensuring the inclusion of all text and image links.

Given the web's vast repository of data, much of it can be irrelevant, potentially obscuring the LLMs' comprehension of key information. This irrelevance can degrade the quality of generated content and lead to unnecessary token consumption. Therefore, it is crucial to filter out non-essential content. Initially, the webpage content is divided into segments based on double newline characters, with each segment treated as an independent chunk. The bge-m3 model is then employed to calculate the similarity score of each chunk relative to the query. A similarity threshold λ is established, and chunks with similarity scores below λ are filtered out, ensuring that only the most relevant information is retained.

To understand the contextual and semantic information of images embedded within the text, we select the contextual content surrounding each `` tag, capturing the relevant text that provides insight into the image's role within the document. We then task vision language models (VLMs) with inferring the caption of the image and describing its content in a concise manner. This process yields two outputs: a descriptive caption and a succinct content summary. These outputs, along with the image's URL, form

¹<https://github.com/infiniflow/ragflow>

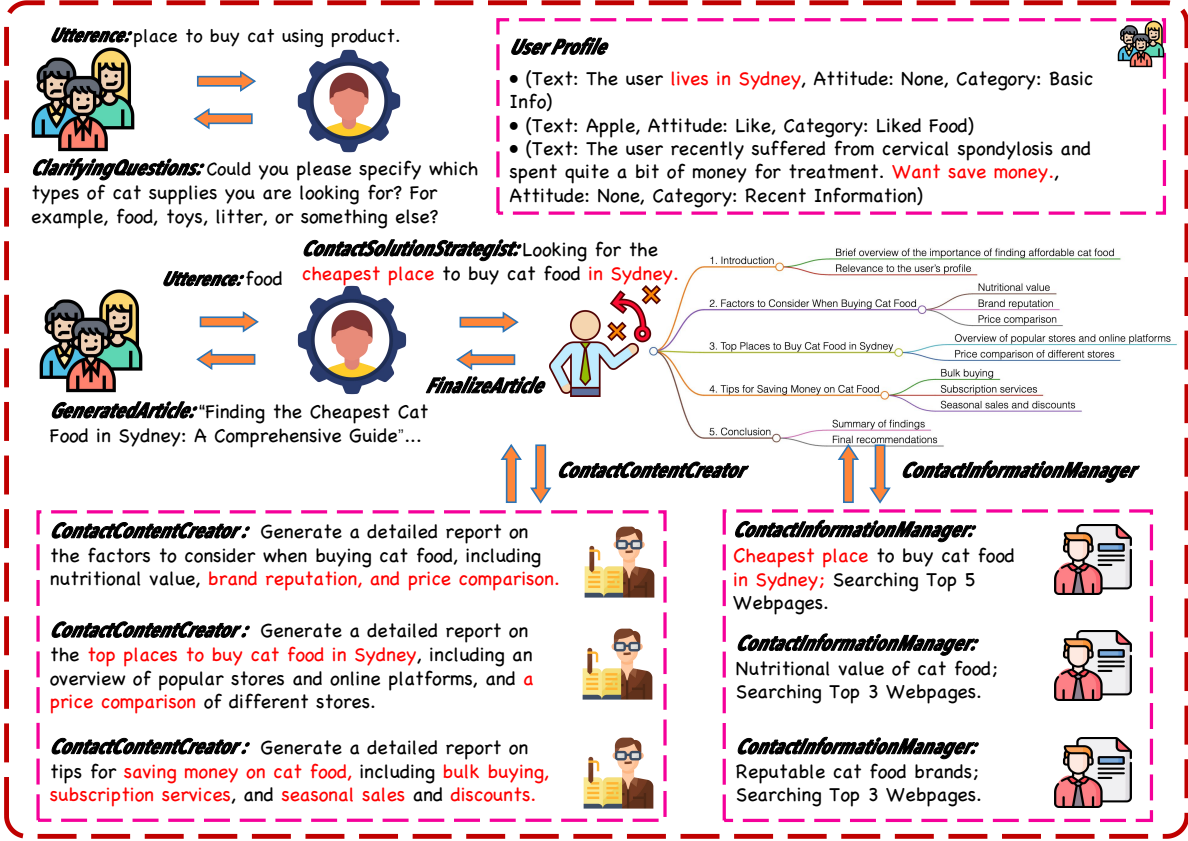


Figure 1: Agent Collaboration Network Framework with a case study "Place to buy cat using product." The red text illustrates how the ACN customizes message passing, the information searching process, and the information generation process to the specific user. The mind map is generated using COT, guiding the function calling for solving the user's information gaining requirements in a logic, deep, and structural way.

a comprehensive set of image-related data. This data is archived for later usage of content generation tasks.

3.4 Content Creator

The Content Creator Agent is responsible for adhering to the creation requirements specified by the Solution Strategist Agent, generating personalized reports with multimodality. This process is meticulously designed to align with the user's profile, ensuring a high degree of personalization and relevance in the generated content.

The LLMs' generation is prompted with user's profile, retrieved external knowledge, and the image information. The profile prompt the generated content aligned with the user's interests, preferences, and needs, ensuring that the generated article is not only informative but also appealing and useful to the user. During the generation, the agent can generate image captions, allowing for the collected image to be integrated smoothly into the report.

3.5 Role Setting

In designing the Account Manager Agent and Solution Strategist, we leveraged the function calling capabilities of LLMs to achieve a higher degree of flexibility. For the Content Creator Agent, which is

tasked solely with text generation, we opted to utilize the traditional text completion ability of LLMs. Detailed function calling settings are available in Appendix ??

3.5.1 Prompt for Account Manager.

Instruction: You are Account Manager in a collaborative agent network aims at providing Personalized Multimodal Information Retrieval and Generation service. Your task is to interact with users in a friendly manner, maintain relationships with customers, ensure customer satisfaction, and understand their needs and expectations through ongoing communication. Furthermore, you are responsible for coordinating the company's Solution Strategist Agent for solving customers' personalized multimedia information retrieval and generation request.

Functions: Normal Reply, Clarifying Questions, Providing Suggestions, Contact Solution Strategist, Tracking User Preferences, Accepting Feedback and Reflection.

3.5.2 Prompt for Solution Strategist.

Instruction: You are Solution Strategist in a collaborative agent network aims at providing Personalized Multimodal Information Retrieval and Generation service. Your task is to develop a logical plan to solve the tasks described in the [User Requirement] conveyed

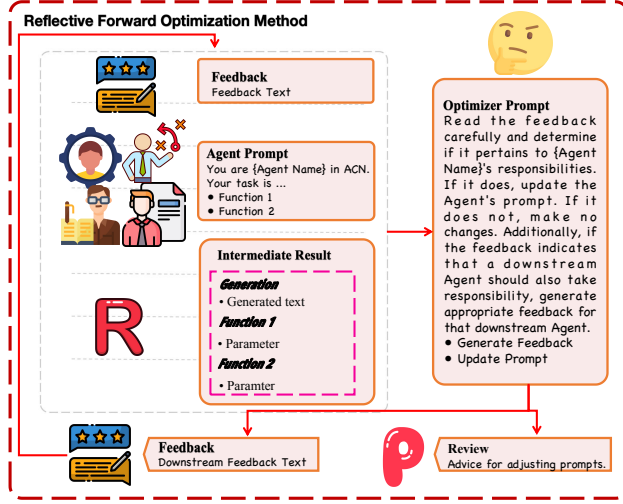


Figure 2: RFO algorithm workflow.

from the Account Manager agent. You should outline this plan step by step, using flexible combinations of calling Search Information and Generate Content. But you must end with the function Finalize Article. Besides, you should also consider the provided *[User Profile]* to make your logical plan specialized for the user.

User Requirement: {User Requirement is here.}

User Profile: {User Profile is here.}

Functions: Search Information, Generate Content, Finalize Article.

3.5.3 Prompt for Content Creator.

Instruction: You are Content Creator Agent in a collaborative agent network aims at providing Personalized Multimodal Information Retrieval and Generation service. Your task is to utilize your professional content creation skills, based on the provided *[External Knowledge]*, and *[Image Source]* as your reading material to generate a detailed multimodal content in markdown format. You should strictly follow the *[Writing Requirement]*, and include appropriate images as much as possible to make the content rich. You also must consider the *[User Profile]*, and makes the content personalized, aligning with user's information, preferences.

External Knowledge: External Knowledge is here.

Image Source: Image Source is here.

Writing Requirement: Writing Requirement is here.

User Profile: User Profile is here.

4 Reflective Forward Optimization

In the process of delivering services to users, they may provide feedback and the Account Manager agent can automatically utilizes to trigger the function of Accepting Feedback and Reflection. This function is essential for the adaptive optimization of the Agent Collaboration Network, enabling real-time, conversational online adjustments. We have developed a novel optimization method termed Reflective Forward Optimization (RFO) to enhance the agent network.

The adjustable parameters for each agent in the ACN are all prompts. Unlike the backpropagation optimization algorithm in neural networks, our RFO is based on a depth-first traversal algorithm for forward propagation optimization. Using an LLM-based

Algorithm 1 Reflective Forward Optimization (RFO)

```

1: Input: User Feedback  $UFB$ , Agent Collaboration Network  $ACN$ ,  $ACN$  Intermediate Result  $RESULT$ 
2: Output: Optimized Agent Collaboration Network  $OACN$ 
3: function RFO( $UFB, ACN, RESULT$ )
4:   Initialize  $FB \leftarrow UFB$ 
5:   Initialize  $stack \leftarrow [(\mathcal{A}, UFB)]$   $\triangleright$  Initialize stack with root agent and user feedback
6:   while  $stack$  is not empty do
7:      $(A, FB) \leftarrow stack.pop()$   $\triangleright$  Current Agent  $A$ , Current Feedback  $FB$ 
8:      $prompt \leftarrow A.prompt$ 
9:      $result \leftarrow RESULT.A$ 
10:     $(Down\_Agents, Down\_FBs, Prompt\_Review) \leftarrow$ 
      Optimizer( $FB, prompt, result$ )
11:     $A.Review\_List \leftarrow Prompt\_Review$ 
12:    for each  $(Down\_Agent, Down\_FB)$  in  $zip(Down\_Agents, Down\_FBs)$  do
13:       $stack.push(Down\_Agent, Down\_FB)$ 
14:    end for
15:  end while
16:  for each Agent in  $ACN$  do
17:    Agent.prompt  $\leftarrow$  UpdatePrompt(Agent.Review_List)
18:  end for
19:  return  $OACN$ 
20: end function

```

optimizer, it systematically reflects on and examines the previous processes used to fulfill user requirements. This allows the algorithm to assign responsibility to agents and make necessary adjustments. The illustrative workflow of RFO is in Figure 2, and the detailed algorithmic process is in Algorithm 1.

The design of optimizer is as follows: **Instruction:** You are an optimizer based on a large language model. The task involves: (1) There is a [Call Agent] that passes parameters [Message] to a [Called Agent]. The external input to the [Call Agent] is [Input], and the output from the [Call Agent] to the external environment is [Output]. (2) The [Call Agent] can adjust the parameters in [Parameter]. (3) The external environment provides [Feedback] on [Output]. Your task is (a) Determine if the cause of the [Feedback] lies with the [Call Agent]. If it does, you need to review each parameter in [Parameter] one by one and provide adjustment suggestions in <review>. (b) If the cause is not with the [Call Agent], you need to provide downstream feedback to the [Called Agent] in the <down_feedback>, to let the [Called Agent] to further reflect himself. If the [Called Agent] is None, then just set <down_feedback> as None. We must use function: Optimize to generate the <review> and <down_feedback>.

Functions: Optimize.

5 Experimental Setup

5.1 Dataset

Personalized dialogue datasets, such as CONVAI2 [9], DuLeMon [26], and KBP [24], emphasize enhancing the personalization of conversations. These datasets typically include a segment of the

user profile, ensuring that the generated dialogues align with this profile. The LaMP [22] benchmark is specifically designed to train and evaluate large language models (LLMs) for personalized outputs, with tasks ranging from personalized title generation to tweet paraphrasing in LaMP4-9. Although these datasets are relevant to our research, they do not fully address our specific focus. Our objective is to design a dataset comprising multiple session chat records, each containing varied-length dialogues between users and AI search engines. Users demonstrate dynamic topic interests, provide feedback, and articulate multimodal information requirements. In response, AI search engines not only engage in basic chat but also generate informative responses.

In light of the lack of appropriate datasets and inspired by recent advancements in utilizing the role-playing capabilities of large language models (LLMs) for simulating realistic scenarios to generate datasets[1, 2], we introduce the synthetic Multi-Session Multi-Turn Personalized Information Inquiry and Generation (MSMTPInfo) dataset. This dataset is meticulously designed to emulate interactions between authentic users, characterized by distinct and evolving personalities, and an AI search engine. Each session comprises multiple turns of dialog between the user and the AI search engine, forming a series of user utterances and corresponding responses.

The dataset spans conversations across 13 diverse main themes and many sub-themes. Initially, we delineate the primary and secondary themes within the conversational content and elucidate the specific actions associated with the attitudes exhibited in user responses, as illustrated in Figure ?? . Following this, we utilize the prompt template shown in Figure ?? to systematically generate data between a user and an AI assistant on a session-by-session basis.

5.2 Baselines and Our Method.

In this study, we evaluate the performance of several notable commercial AI search engines. Besides, given the impracticality of conducting fair comparisons due to discrepancies in the number of indexed webpages, variations in backbone LLM models, and the inability to modify these closed-source commercial engines for rigorous studies, we also consider one open-source AI search engine.

1. Perplexity²: Widely-used and popular all over the world.
2. TianGong³: It has research mode capabilities that produce meticulously logical and comprehensive reports.
3. Perplexica⁴: It can be considered as open-source counterpart of Perplexity. We set each query searching 5 webpages via Bing, utilizing GPT-4o-mini as the backbone LLM.
4. Our ACN: we integrated our proposed ACN framework into the open-soure Perplexica framework, supporting multimodal text and image outputs, personalized information generation, logical planning for complex queries, and adaptive online adjustments.

5.3 Evaluation Setup.

Our evaluation methodology leverages insights from prior research on using LLM-as-a-judge [31] method, which demonstrated that employing GPT-4 as an adjudicator aligns with human assessments at an agreement rate exceeding 80%. In our approach, we present

the LLM-based judger with two responses: one generated by the ACN and the other from alternative configurations. The judger is tasked with distinguishing between these responses and delivering a judgment based on predefined criteria of win, loss, or tie. To mitigate any positional bias inherent in LLMs, we subsequently reverse the positions of the two responses within the prompt and re-evaluate. If the judger’s decisions are consistent across both evaluations—whether both are wins, ties, or losses—the final judgment is accordingly recorded as a win, tie, or loss. Conversely, if the adjudications differ, the final result is deemed a tie. Finally, we calculate the win rate, tie rate, loss rate of ACN compared to other alternative configuration. Our evaluation encompasses multiple dimensions to provide a comprehensive assessment:

- **Content Richness**: The richness of the content, including both textual and visual elements, is vital for capturing and sustaining user interest.
- **Information Usefulness**: Despite the richness of content, the actual utility of the information is paramount. An abundance of redundant or irrelevant information fails to meet the user’s need for efficient knowledge acquisition. Hence, the AI search engine’s effectiveness should be further evaluated from the perspective of information usefulness.
- **Content Personalization**: Leveraging user profile information revealed either in prior sessions or during the ongoing interaction is essential for customizing the dialogue. This personalized approach is a key factor in delivering a superior user experience through the AI search engine.
- **Writing Logicality**: AI search engines can address complex queries that traditional search engines fail to resolve. These questions often require more than simple keyword searches, as they encompass intricate problems with inherent logical structures.

6 Experimental Results

We now present our part of experimental results, and report findings from various auxiliary studies and analyses.

6.1 Richness Analysis

Our proposed ACN demonstrates a substantial superiority over current AI search engines that are limited to text-based article generation. By incorporating multiple modalities, including text, appropriate image insertions, and tabular data presentation, ACN-generated articles offer enhanced visual appeal and engagement. This multimodal approach not only captivates readers’ attention but also provides a better satisfying experience. Figure ?? exemplifies this advantage.

6.2 Usefulness and Personalization Analysis

The whole dataset’s comparative assessment outcomes in different topics are statistically summarized and visualized via a radar chart in Figure 4. The left subplot of the figure exclusively concentrates on criteria usefulness. Subsequently, right subplot of the figure evaluates personalization by assessing how well the responses aligned with the user’s profile. It was observed that response A and response B often reached a draw across various topics, with instances of unilateral victories or defeats being comparatively rare. However, when taking into personalization, response B consistently

²<https://www.perplexity.ai/>

³<https://www.tiangong.cn/>

⁴<https://github.com/ItzCrazyKns/Perplexica/tree/v1.3.0>

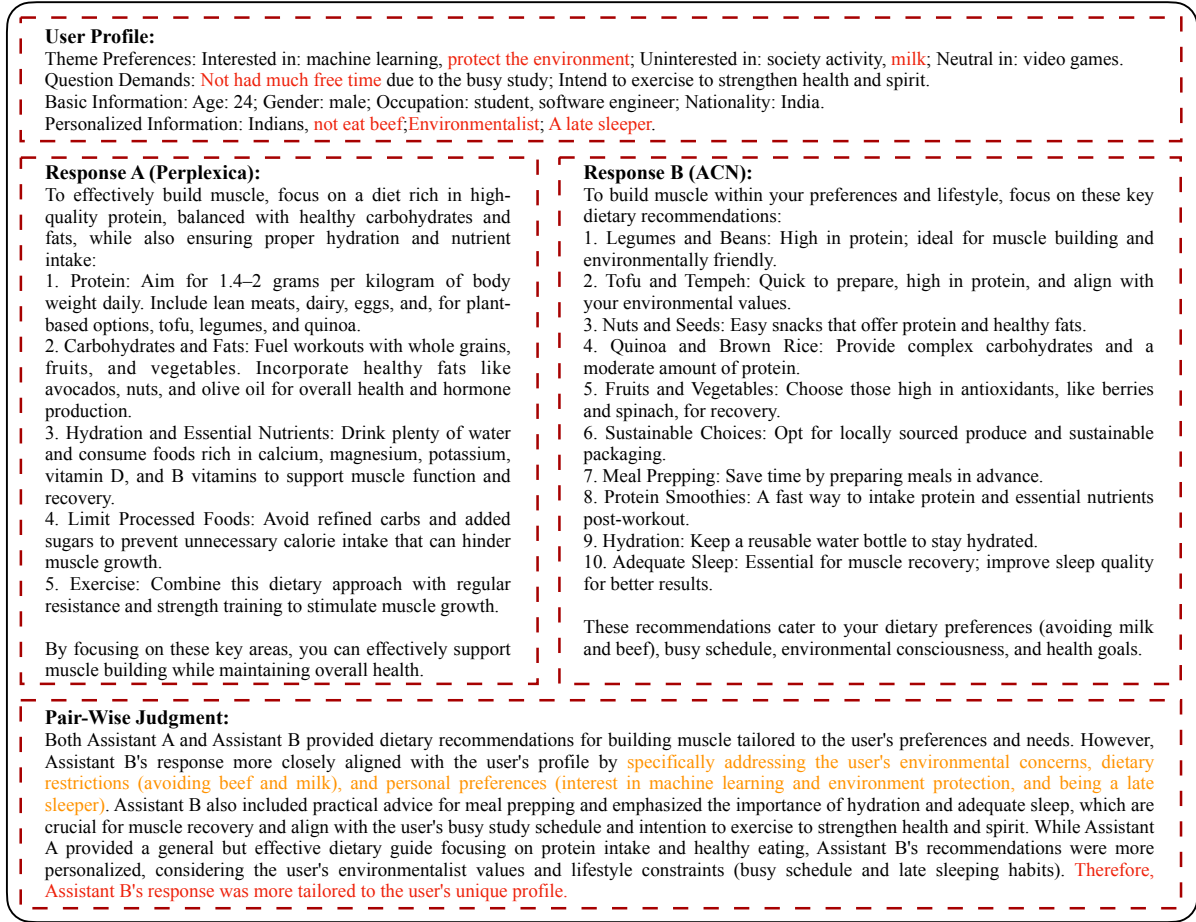


Figure 3: Comparison of AI Search Engine Responses to the Query "Give me a dietary recommendation for building muscle." A LLM played judge subsequently determines that Response B (ACN) is better.

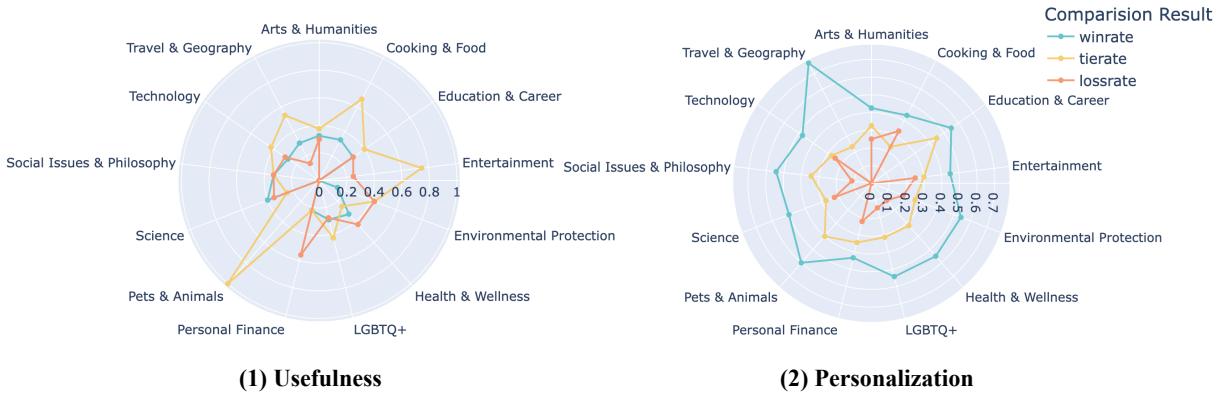


Figure 4: The results of pairwise comparisons between Basic and ACN responses across all categories on the MSMTPIInfo dataset.

outperformed response A, demonstrating a superior ability of ACN delivering tailored responses that align more closely with the user's unique profile.

Further, a case study of the pair-wise judgement result is illustrated in Figure 3. This case study provides a brief summary of two AI search engines' responses to the identical user query. The

response generated by Perplexica (anonymized as Response A) does not take the user profile into account. In contrast, the response generated by the ACN (anonymized as Response B) consistently tracks the user profile and produces a more customized recommendation. This investigation elucidates the comparative analysis of response between two AI search engines. The focal point of this case study

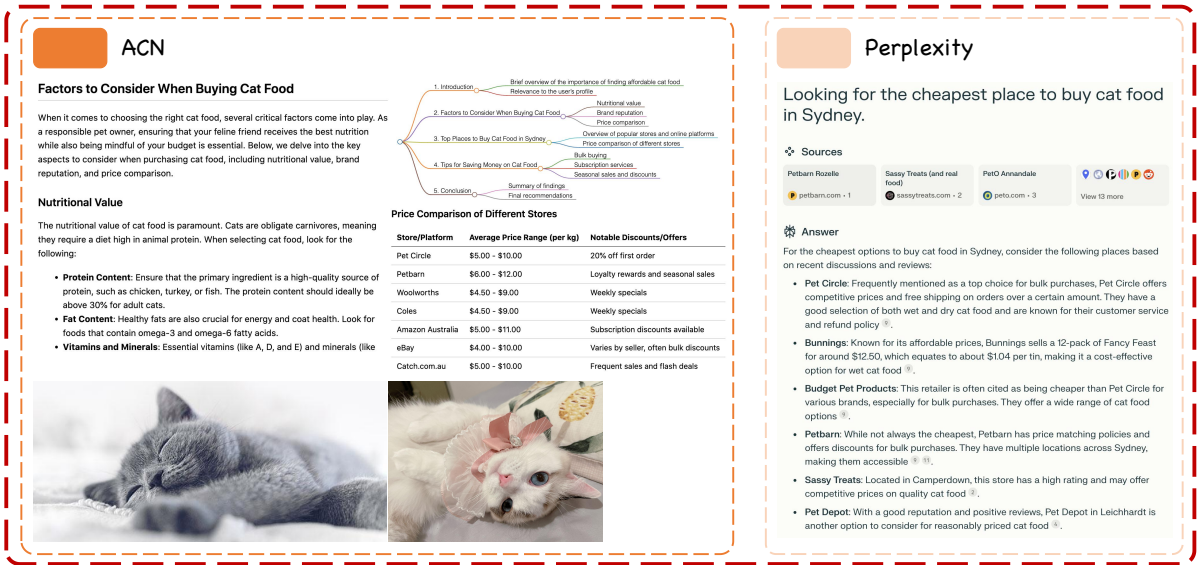


Figure 5: Comparison between ACN and Perplexity.

hinges on the assessment of responses provided by A and B in the context of their consideration of the user profile. The judgement result unveils that while both responses delivered accurate and comprehensive response. But response B is distinguished by the personalized nature, tailored specifically to the user’s preferences and needs.

6.3 Logicality Analysis

We have developed a Solution Strategist Agent that leverages Chain of Thought (COT) reasoning to enhance the logical capabilities of search engines in generating responses to complex questions. To rigorously evaluate the logical soundness of ACN, we assess the following three dimensions:

- **Depth:** This metric evaluates the thoroughness of the strategic plan in exploring a particular point of consideration.
- **Comprehensiveness:** This metric measures the extent to which the strategic plan addresses both explicit and implicit factors necessary for a robust response.
- **Reasonability:** This metric assesses the relative rationality of the strategic plan when analogized to human problem-solving and planning processes.

We benchmarked our ACN against TianGong and Perplexity. The pair-wise comparison results are illustrated in Figure 6, demonstrate a significant absolute improvement, underscoring the critical role of the Solution Strategist Agent and affirming the superiority of our proposed ACN.

7 Discussion and Future Works

Our current research has yet to undergo rigorous empirical validation, particularly concerning the ACN’s capability of online learning and prompt adjusting based on user feedback. This limitation arises from the necessity for real human evaluations, which require more time and carefully designed experimental protocols. To address this gap, our future research will focus on the following:

1. Identifying suitable volunteers for testing the ACN.

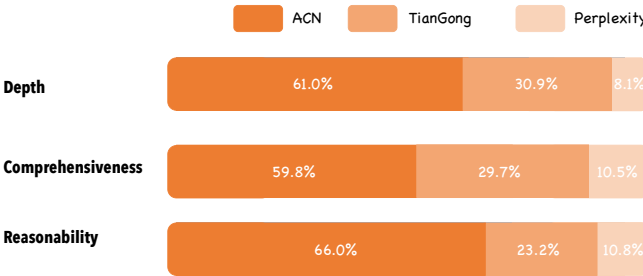


Figure 6: Comparative Evaluation Results of Pairwise Comparisons among ACN, TianGong, and Perplex. We conducted pairwise comparisons among these search engines, calculating the adjusted win rate for each and subsequently normalizing the results.

2. Standardizing experimental procedures by designing feedback types. Users will be able to provide feedback within predefined categories.
3. Comparing the ACN’s performance with other models such as TianGong and Perplexity. Given that large models possess context-aware prompt learning capabilities, they theoretically offer some degree of real-time adjustment. However, the extent of this capability remains unclear and necessitates empirical investigation.
4. Developing metrics to evaluate the AI search engine’s responsiveness to user feedback. We suppose that feedback will influence dialogue consistency and content personalization.

Acknowledgments

This work was sponsored by the Australian Research Council under the Linkage Projects Grant LP210100129, and by the program of China Scholarships Council (No. 202308200014).

References

- [1] Zahra Abbasiantaeb, Yifei Yuan, Evangelos Kanoulas, and Mohammad Alian-nejadi. 2024. Let the LLMs Talk: Simulating Human-to-Human Conversational QA via Zero-Shot LLM-to-LLM Interactions. In *Proceedings of the 17th ACM International Conference on Web Search and Data Mining (Merida, Mexico) (WSDM '24)*. Association for Computing Machinery, New York, NY, USA, 8–17. <https://doi.org/10.1145/3616855.3635856>
- [2] Yelaman Abdullin, Diego Molla-Aliod, Bahadorreza Ofoghi, John Yearwood, and Qingyang Li. 2024. Synthetic Dialogue Dataset Generation using LLM Agents. In *Proceedings of the GEM Workshop at EMNLP 2023*. arXiv:2401.17461 [cs.CL] <https://arxiv.org/abs/2401.17461>
- [3] Jinheon Baek, Nirupama Chandrasekaran, Silviu Cucerzan, Allen Herring, and Sujay Kumar Jauhar. 2024. Knowledge-augmented large language models for personalized contextual query suggestion. In *Proceedings of the ACM on Web Conference 2024*. 3355–3366.
- [4] Sebastian Borgeaud, Arthur Mensch, Jordan Hoffmann, Trevor Cai, Eliza Rutherford, Katie Millican, George Bm Van Den Driessche, Jean-Baptiste Lespiau, Bogdan Damoc, Aidan Clark, et al. 2022. Improving language models by retrieving from trillions of tokens. In *International conference on machine learning*. PMLR, 2206–2240.
- [5] Samuel R. Bowman, Gabor Angeli, Christopher Potts, and Christopher D. Manning. 2015. A large annotated corpus for learning natural language inference. In *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing*, Lluís Màrquez, Chris Callison-Burch, and Jian Su (Eds.). Association for Computational Linguistics, Lisbon, Portugal, 632–642. <https://doi.org/10.18653/v1/D15-1075>
- [6] Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. Language Models are Few-Shot Learners. arXiv:2005.14165 [cs.CL]
- [7] Xin Cheng, Di Luo, Xiuying Chen, Lemao Liu, Dongyan Zhao, and Rui Yan. 2024. Lift yourself up: Retrieval-augmented text generation with self-memory. *Advances in Neural Information Processing Systems* 36 (2024).
- [8] Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, Parker Schuh, Kensen Shi, Sasha Tsvyashchenko, Joshua Maynez, Abhishek Rao, Parker Barnes, Yi Tay, Noam Shazeer, Vinodkumar Prabhakaran, Emily Reif, Nan Du, Ben Hutchinson, Reiner Pope, James Bradbury, Jacob Austin, Michael Isard, Guy Gur-Ari, Pengcheng Yin, Toju Duke, Anselm Levskaya, Sanjay Ghemawat, Sunipa Dev, Henryk Michalewski, Xavier Garcia, Vedant Misra, Kevin Robinson, Liam Fedus, Denny Zhou, Daphne Ippolito, David Luan, Hyeontaek Lim, Barret Zoph, Alexander Spiridonov, Ryan Sepassi, David Dohan, Shivani Agrawal, Mark Omernick, Andrew M. Dai, Thanumalayan Sankaranarayanan Pillai, Marie Pellat, Aitor Lewkowycz, Erica Moreira, Rewon Child, Oleksandr Polozov, Katherine Lee, Zongwei Zhou, Xuezhi Wang, Brennan Saeta, Mark Diaz, Orhan Firat, Michele Catasta, Jason Wei, Kathy Meier-Hellstern, Douglas Eck, Jeff Dean, Slav Petrov, and Noah Fiedel. 2022. PaLM: Scaling Language Modeling with Pathways. arXiv:2204.02311 [cs.CL]
- [9] Emily Dinan, Varvara Logacheva, Valentin Malykh, Alexander Miller, Kurt Shuster, Jack Urbanek, Douwe Kiela, Arthur Szlam, Iulian Serban, Ryan Lowe, Shrimai Prabhumoye, Alan W Black, Alexander Rudnicky, Jason Williams, Joelle Pineau, Mikhail Burtsev, and Jason Weston. 2019. The Second Conversational Intelligence Challenge (ConvAI2). arXiv:1902.00098 [cs.AI] <https://arxiv.org/abs/1902.00098>
- [10] Luyu Gao, Zhuoyun Dai, Panupong Pasupat, Anthony Chen, Arun Tejasvi Chaganty, Yicheng Fan, Vincent Zhao, Ni Lao, Hongrae Lee, Da-Cheng Juan, and Kelvin Guu. 2023. RARR: Researching and Revising What Language Models Say, Using Language Models. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, Anna Rogers, Jordan Boyd-Graber, and Naoaki Okazaki (Eds.). Association for Computational Linguistics, Toronto, Canada, 16477–16508. <https://doi.org/10.18653/v1/2023.acl-long.910>
- [11] Shen Gao, Yuntao Wen, Minghang Zhu, Jianing Wei, Yuhang Cheng, Qunzi Zhang, and Shuo Shang. 2024. Simulating Financial Market via Large Language Model based Agents. arXiv:2406.19966 [cs.CL] <https://arxiv.org/abs/2406.19966>
- [12] Yunfan Gao, Yun Xiong, Xinyu Gao, Kangxiang Jia, Jinliu Pan, Yuxi Bi, Yi Dai, Jiawei Sun, Qianyu Guo, Meng Wang, and Haofen Wang. 2024. Retrieval-Augmented Generation for Large Language Models: A Survey. arXiv:2312.10997 [cs.CL]
- [13] Wei Ji, Hao Fei, Yinwei Wei, Zhedong Zheng, Juncheng Li, Long Chen, Lizi Liao, Yueting Zhuang, and Roger Zimmermann. 2024. MMGR'24: 2024 Workshop on Deep Multimodal Generation and Retrieval. In *Proceedings of the 32th ACM International Conference on Multimedia Workshop*.
- [14] Gangwoo Kim, Sungdong Kim, Byeongguk Jeon, Joonsuk Park, and Jaewoo Kang. 2023. Tree of Clarifications: Answering Ambiguous Questions with Retrieval-Augmented Large Language Models. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, Houda Bouamor, Juan Pino, and Kalika Bali (Eds.). Association for Computational Linguistics, Singapore, 996–1009. <https://doi.org/10.18653/v1/2023.emnlp-main.63>
- [15] Hannah Rose Kirk, Bertie Vidgen, Paul Röttger, and Scott A. Hale. 2023. Personalisation within bounds: A risk taxonomy and policy framework for the alignment of large language models with personalised feedback. arXiv:2303.05453 [cs.CL] <https://arxiv.org/abs/2303.05453>
- [16] Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, et al. 2020. Retrieval-augmented generation for knowledge-intensive nlp tasks. *Advances in Neural Information Processing Systems* 33 (2020), 9459–9474.
- [17] Binxu Li, Tiankai Xu, Yuanting Pan, Zhe Xu, Jie Luo, Ruiyang Ji, Shilong Liu, Haoyu Dong, Zihao Lin, and Yixin Wang. 2024. MMedAgent: Learning to Use Medical Tools with Multi-modal Agent. arXiv:2407.02483 [cs.CL] <https://arxiv.org/abs/2407.02483>
- [18] Wenhan Liu, Yujia Zhou, Yutao Zhu, and Zhicheng Dou. 2024. How to personalize and whether to personalize? Candidate documents decide. *Knowledge and Information Systems* (2024), 1–24.
- [19] Xinbei Ma, Yeyun Gong, Pengcheng He, Hai Zhao, and Nan Duan. 2023. Query Rewriting for Retrieval-Augmented Large Language Models. arXiv:2305.14283 [cs.CL]
- [20] Yujia Qin, Shihao Liang, Yining Ye, Kunlun Zhu, Lan Yan, Yaxi Lu, Yankai Lin, Xin Cong, Xiangru Tang, Bill Qian, et al. 2023. Toolllm: Facilitating large language models to master 16000+ real-world apis. arXiv preprint arXiv:2307.16789 (2023).
- [21] Chris Richardson, Yao Zhang, Kellen Gillespie, Sudipta Kar, Arshdeep Singh, Zeynab Raeesy, Omar Zia Khan, and A. Sethy. 2023. Integrating Summarization and Retrieval for Enhanced Personalization via Large Language Models. *ArXiv abs/2310.20081* (2023). <https://doi.org/10.48550/arXiv.2310.20081>
- [22] Alireza Salemi, Sheshera Mysore, Michael Bendersky, and Hamed Zamani. 2024. LaMP: When Large Language Models Meet Personalization. arXiv:2304.11406 [cs.CL] <https://arxiv.org/abs/2304.11406>
- [23] Yunxiao Shi, Xing Zi, Zijing Shi, Haimin Zhang, Qiang Wu, and Min Xu. 2024. Enhancing Retrieval and Managing Retrieval: A Four-Module Synergy for Improved Quality and Efficiency in RAG Systems. arXiv:2407.10670 [cs.CL] <https://arxiv.org/abs/2407.10670>
- [24] Hongru Wang, Minda Hu, Yang Deng, Rui Wang, Fei Mi, Weichao Wang, Yasheng Wang, Wai-Chung Kwan, Irwin King, and Kam-Fai Wong. 2023. Large Language Models as Source Planner for Personalized Knowledge-grounded Dialogue. arXiv:2310.08840 [cs.CL]
- [25] Shuting Wang, Zhicheng Dou, Jing Yao, Yujia Zhou, and Ji-Rong Wen. 2023. Incorporating Explicit Subtopics in Personalized Search. In *Proceedings of the ACM Web Conference 2023* (Austin, TX, USA) (WWW '23). Association for Computing Machinery, New York, NY, USA, 3364–3374. <https://doi.org/10.1145/3543507.3583488>
- [26] Xinchao Xu, Zhibin Gou, Wenquan Wu, Zheng-Yu Niu, Hua Wu, Haifeng Wang, and Shihang Wang. 2022. Long Time No See! Open-Domain Conversation with Long-Term Persona Memory. arXiv preprint arXiv:2203.05797 (2022).
- [27] Rui Yang, Lin Song, Yanwei Li, Sijie Zhao, Yixiao Ge, Xiu Li, and Ying Shan. 2024. Gpt4tools: Teaching large language model to use tools via self-instruction. *Advances in Neural Information Processing Systems* 36 (2024).
- [28] Ori Yoran, Tomer Wolfson, Ori Ram, and Jonathan Berant. 2023. Making Retrieval-Augmented Language Models Robust to Irrelevant Context. arXiv:2310.01558 [cs.CL]
- [29] Zihan Zhang, Meng Fang, and Ling Chen. 2024. RetrievalQA: Assessing Adaptive Retrieval-Augmented Generation for Short-form Open-Domain Question Answering. arXiv:2402.16457 [cs.CL] <https://arxiv.org/abs/2402.16457>
- [30] Andrew Zhao, Daniel Huang, Quentin Xu, Matthieu Lin, Yong-Jin Liu, and Gao Huang. 2023. Expel: Llm agents are experiential learners. In *Proceedings of THE 37th Association for the Advancement of Artificial Intelligence Conference*.
- [31] 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 (2024).
- [32] Wangchunshu Zhou, Yixin Ou, Shengwei Ding, Long Li, Jialong Wu, Tiannan Wang, Jiamin Chen, Shuai Wang, Xiaohua Xu, Ningyu Zhang, Huajun Chen, and Yuchen Eleanor Jiang. 2024. Symbolic Learning Enables Self-Evolving Agents. arXiv:2406.18532 [cs.CL] <https://arxiv.org/abs/2406.18532>
- [33] Yujia Zhou, Zhicheng Dou, and Ji-Rong Wen. 2020. Encoding history with context-aware representation learning for personalized search. In *Proceedings of the 43rd international ACM SIGIR conference on research and development in information retrieval*. 1111–1120.
- [34] Yujia Zhou, Qiannan Zhu, Jiajie Jin, and Zhicheng Dou. 2024. Cognitive Personalized Search Integrating Large Language Models with an Efficient Memory Mechanism. In *Proceedings of the ACM on Web Conference 2024* (Singapore, Singapore) (WWW '24). Association for Computing Machinery, New York, NY, USA, 1464–1473. <https://doi.org/10.1145/3589334.3645482>