

ResearchAgent: Iterative Research Idea Generation over Scientific Literature with Large Language Models

Anonymous ACL submission

Abstract

Scientific Research, vital for improving human life, is hindered by its inherent complexity, slow pace, and the need for specialized experts. To enhance its productivity, we propose a ResearchAgent, a Large Language Model (LLM)-powered research idea generation agent, which automatically defines problems, proposes methods and designs experiments, while iteratively refining them based on the feedback from LLM-powered reviewing agents. Specifically, starting with a core scientific paper as the primary focus to generate ideas, our ResearchAgent is augmented not only with relevant publications by connecting information over an academic graph but also entities retrieved from an entity-centric knowledge store based on their shared underlying concepts, mined across numerous papers. Then, mimicking the human approach to iteratively improving ideas with peer discussions, we leverage multiple ReviewingAgents that provide reviews and feedback via iterative revision processes. These reviewing agents are instantiated with human preference-aligned LLMs whose criteria for evaluation are elicited from actual human judgments via LLM prompting. We experimentally validate our ResearchAgent on scientific publications across multiple disciplines, showing its effectiveness in generating novel, clear, and valid ideas based on human and model-based evaluation results.

1 Introduction

Scientific research plays a crucial role in driving innovation, advancing knowledge, solving problems, expanding our understanding of the world, and ultimately improving the lives of people in tangible ways. This process usually consists of two key components: the formulation of new research ideas and the validation of these ideas through well-crafted experiments, which are typically conducted by human researchers (Hope et al., 2023; Wang et al., 2023a; Huang et al., 2023). However, this is a

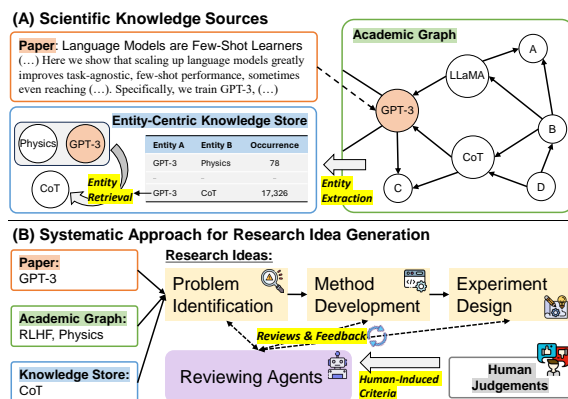


Figure 1: (A) The scientific knowledge used for research idea generation consists of a paper, its relationships over an academic graph, and entities within a knowledge store extracted from numerous papers. (B) Given them, the proposed research idea generation process involves problem identification, method development, and experiment design. Those are also iteratively refined by reviews and feedback from reviewing agents, aligned with criteria induced from human judgements.

tedious process, which requires reading and synthesizing overwhelming amounts of knowledge over the vast corpus of rapidly growing scientific literature to formulate research ideas, but also designing and performing experimental validations of those ideas. For example, the number of academic papers published per year is more than 7 million (Fire and Guestrin, 2019). Also, the process of testing a new pharmaceutical drug is labor-intensive, often taking several years (Vamathevan et al., 2019). These constraints highlight the potential benefits of integrating AI assistance to enhance the efficiency and productivity of scientific research.

Recently, Large Language Models (LLMs) (Touvron et al., 2023; OpenAI, 2023; Anil et al., 2023) have shown impressive capabilities in processing and generating text with remarkable accuracy, even outperforming human experts across diverse specialized domains including math, physics, history, law, medicine, and ethics. Thus, LLMs may be a transformative tool to accelerate the scientific research process, helping humans perform it. Specifically, LLMs can process and analyze large volumes

of data at a speed and scale that far exceeds human capabilities, but also identify patterns, trends, and correlations that may not be immediately apparent to human researchers. This may enable them to identify novel research opportunities that might otherwise remain undiscovered. Moreover, LLMs can assist in experimental validation by conducting the experiments and interpreting the results, thereby significantly accelerating the research cycle. In this work, our focus is on the first phase of scientific research, namely *research idea generation*, which involves problem identification, method development, and experiment design.

While there are few recent works in the domain of LLM-augmented scientific discovery, they focus on largely different scenarios. Specifically, most of them (Huang et al., 2023; AI4Science and Quantum, 2023; Bran et al., 2023) have mainly targeted accelerating the experimental validation processes (phase 2 of the scientific research), by writing the code for machine-learning models, facilitating the exploration of chemical spaces, or advancing the simulation of molecular dynamics. On the other hand, the usage of LLMs in the initial phase of research idea generation, whose key focus is on conceptualizing new scientific questions, methodologies, and experiments, remains underexplored. We note that, along this line of work, few recent methods (Wang et al., 2023b; Yang et al., 2023; Qi et al., 2023) have studied the problem of hypothesis generation, which is based on Literature-based Discovery (LBD) (Swanson, 1986). However, this setting is suboptimal and not fully open-ended, restricted to forming new relationships between two concepts (which are sometimes defined in natural language format), like the potential applications of a new medication for a specific ailment. Also, its scope is narrow, lacking consideration of the wider processes involved in scientific idea generation.

In this work, we aim to build an LLM-powered research agent, which is capable of generating research ideas over scientific literature. Specifically, mirroring the human approach to formulating research ideas, the proposed agent begins by reading an academic paper, then explores related papers based on references and citation relationships. However, despite its simplicity and straightforwardness, the fact that the agent focuses only on one paper and its immediate references could hinder its ability to fully grasp and utilize the broader contextual knowledge of relevant scientific fields. It is worth noting that this contextual knowledge

is accumulated over numerous papers (oftentimes across multiple disciplines), which skilled human researchers either possess, imbibe through communication with other researchers, or learn from perusals of scientific literature, then leverage to come up with and develop new ideas. In addition, another limitation of this one-step generation approach (that concludes once the ideas are formulated) is the lack of an iterative refinement process based on reviews and feedback from multiple perspectives, which differs from typical human-driven research processes, which develop and improve research ideas through multiple peer discussions.

To tackle those limitations, we further propose to expand the idea generation process by not only augmenting it with the knowledge retrieved from an entity-centric knowledge store but also iteratively refining the generated ideas through collaborative efforts with LLM-powered multiple reviewing agents. More specifically, we first construct a knowledge store, which finds and aggregates entity co-occurrences from scientific articles. This entity-centric knowledge store thus captures the mutual relevance between different entities, and is used for retrieving the knowledge that is not present within the accessed articles but may be relevant to them through underlying concepts and principles, which provides valuable insights for our idea generation. Also, to enhance generated research ideas with iterative improvements, we design multiple reviewing agents (based on LLMs), each of which generates a review and feedback on the developed ideas, with their own evaluation criteria. We note that those evaluation criteria are induced by human judgments, to align the LLM-based automatic evaluations with actual human preferences. Then, based on the generated reviews and feedback, the proposed LLM-powered research agent is prompted again to refine the areas for improvement. We refer to our overall framework as ResearchAgent, which is illustrated in Figure 1.

We experimentally validate the effectiveness of ResearchAgent for research idea generation based on scientific literature across multiple disciplines. Then, on a battery of tests conducted with human- and model-based evaluations, ResearchAgent outperforms strong LLM-powered baselines by large margins, generating more clear, relevant, and significant ideas that are especially novel. Moreover, further analyses demonstrate the efficacy of augmenting ResearchAgent with the entity-centric knowledge store and the iterative idea refinement steps.

2 Related Work

Large Language Models Large Language Models (LLMs), which are trained on massive text corpora with multi-billion parameters and through various training strategies (such as pre-training, fine-tuning, and reinforcement learning), have shown impressive performances across a wide range of tasks (OpenAI, 2023; Anil et al., 2023). Their capability extends to advanced scientific fields, which include mathematics, physics, medicine, and computer science (Romera-Paredes et al., 2023; Bran et al., 2023; Huang et al., 2023). A recent study on GPT-4 shows that it is capable of understanding DNA sequences, designing biomolecules, predicting the behavior of molecular systems, and solving Partial Differential Equation (PDE) problems (AI4Science and Quantum, 2023). However, they have mainly been used for accelerating the experimental validation of already identified research hypotheses, but not for identifying new problems.

Hypothesis Generation The principle of hypothesis generation is based on literature-based discovery (Swanson, 1986), which aims to discover relationships between concepts (Henry and McInnes, 2017). For instance, these concepts could be a specific disease and a compound not yet considered as a treatment for it. Early works on automatic hypothesis generation first build a corpus of discrete concepts, and then identify their relationships with machine learning approaches, e.g., using similarities between word (concept) vectors (Tshitoyan et al., 2019) or applying link prediction methods over a graph (where concepts are nodes) (Sybrandt et al., 2020; Nadkarni et al., 2021). Recent approaches are further powered by LLMs (Wang et al., 2023b; Qi et al., 2023; Yang et al., 2023), leveraging their prior knowledge about scientific disciplines. However, the aforementioned approaches are designed to identify potential relationships between two variables or generate sentences of two relations, which may be sub-optimal to capture the complexity and multifaceted nature of real-world problems. Yet, we target challenging and more open-ended scenarios, aiming to generate research ideas that involves comprehensive processes of formulating problems, methods, and experiment designs. Also, during generation, our approach leverages a store of accumulated knowledge extracted from vast amounts of scientific literature, which goes beyond prior work that is based on singular, non-accumulated data (such as the individual abstracts of cited papers).

Knowledge-Augmented LLMs The approach to augment LLMs with external knowledge enhances their utility, making them more accurate and relevant to specific target contexts. Much prior work aims at improving the factuality of LLM responses to given queries by retrieving the relevant documents and then injecting them into the input of LLMs (Lazaridou et al., 2022; Ram et al., 2023; Shi et al., 2023). In addition, given that entities or facts are atomic units for representing knowledge, recent studies further augment LLMs with them (Baek et al., 2023; Wu et al., 2023). In contrast to these efforts which use knowledge units piecemeal, we instead jointly leverage accumulated knowledge over massive troves of scientific papers. More recently, Baek et al. (2024) proposes to use accumulated entities (extracted from various web search contexts) for query suggestion, which yet has a different objective that aims to narrow the focus of LLMs to entities already present in the given context of LLMs. Instead, our approach retrieves and integrates entities outside the given context yet relevant to it, enabling LLMs to explore other concepts or subjects for fruitful idea generation.

Iterative Refinements with LLMs Similar to humans, LLMs do not always generate the optimal outputs on their first attempt, but humans can iteratively refine what they generate through feedback from themselves and others. Motivated by this, a large volume of recent studies (including hypothesis generation work) have investigated the potential of LLMs to correct and refine their outputs, showing they indeed possess those capabilities (Welleck et al., 2023; Madaan et al., 2023; Shridhar et al., 2023; Ganguli et al., 2023; Wang et al., 2023b; Qi et al., 2023; Yang et al., 2023). Based on their findings, we extend this (and further test its capability) to our novel scenario of research problem, method, and experiment design generation processes.

3 Method

We present ResearchAgent, a system that automatically proposes research ideas with LLMs.

3.1 LLM-Powered Research Idea Generation

We begin with formally introducing the new problem of our research idea generation, followed by explaining LLMs used as a basis to tackle it.

Research Idea Generation The goal of the research idea generation task is to formulate new and valid research ideas, to enhance the overall efficiency of the first phase of scientific discovery,

which consists of three systematic steps: identifying problems, developing methods, and designing experiments. We note that this three-step process mirrors human research practices, capturing our approach to exploring new problems, crafting innovative solutions, and testing our ideas, constituting a cycle of questioning, innovating, and validating. Specifically, we first identify problems by noting gaps or contradictions in current knowledge. Following problem identification, we devise methodologies using relevant procedures and tools. The final stage involves experiment design, setting up tests to validate our hypotheses.

To accomplish the aforementioned steps, the existing literature (e.g., academic publications) is used as a primary source, which provides insights about existing knowledge along with gaps and unanswered questions. Formally, let \mathcal{L} be the literature, and \mathbf{o} be the ideas that consist of the problem \mathbf{p} , method \mathbf{m} , and experiment design \mathbf{d} , as follows: $\mathbf{o} = [\mathbf{p}, \mathbf{m}, \mathbf{d}]$ where each item consists of a sequence of tokens and $[\cdot]$ denotes a concatenation operation. Then, the idea generation model f can be represented as follows: $\mathbf{o} = f(\mathcal{L})$, which is further decomposed into three submodular steps: $\mathbf{p} = f(\mathcal{L})$ for identifying problems, $\mathbf{m} = f(\mathbf{p}, \mathcal{L})$ for developing methods, and $\mathbf{d} = f(\mathbf{p}, \mathbf{m}, \mathcal{L})$ for designing experiments. In this work, we operationalize f with LLMs, leveraging their capability to understand and generate academic text.

Large Language Models Before describing the LLM in the context of our problem setup, let us first provide its general definition, which takes an input sequence of tokens \mathbf{x} and generates an output sequence of tokens \mathbf{y} , as follows: $\mathbf{y} = \text{LLM}_\theta(\mathcal{T}(\mathbf{x}))$. Here, the model parameters θ are typically fixed after training, due to the high costs of further fine-tuning. In addition, the prompt template \mathcal{T} serves as a structured format that outlines the context (including the task descriptions and instructions) to direct the model in generating the desired outputs.

3.2 Knowledge-Augmented LLMs for Research Idea Generation

We now turn to our primary focus of automatically generating research ideas with LLMs. Recall that we aim to produce a complete idea consisting of the problem, method, and experiment design ($\mathbf{o} = [\mathbf{p}, \mathbf{m}, \mathbf{d}]$), while using the existing literature \mathcal{L} as a primary source of information. We operationalize this with LLMs by instantiating the aforementioned research idea genera-

tion function f with LLM coupled with the task-specific template. Formally, $\mathbf{p} = \text{LLM}(\mathcal{T}_p(\mathcal{L}))$ indicates the problem identification step, followed by $\mathbf{m} = \text{LLM}(\mathcal{T}_m(\mathbf{p}, \mathcal{L}))$ for method development and $\mathbf{d} = \text{LLM}(\mathcal{T}_e(\mathbf{p}, \mathbf{m}, \mathcal{L}))$ for experiment design, which constitutes the full idea: $\mathbf{o} = [\mathbf{p}, \mathbf{m}, \mathbf{d}]$.

Following this general formulation, the important question to answer is how is the massive literature used for actually generating the research ideas with LLMs. It is worth noting that, due to the constraints of their input lengths and their reasoning abilities, particularly over long contexts (Liu et al., 2023), it is not possible to incorporate all the existing publications from the literature \mathcal{L} into the LLM input. Instead, we should find a meaningful subset from them. To achieve this, we mirror the process followed by human researchers, who expand their knowledge of a paper by perusing other papers that either cite or are cited by it. Similarly, for LMM, we initiate its literature review process by providing a core paper l_0 from \mathcal{L} and then selectively incorporating subsequent papers $\{l_1, \dots, l_n\}$ that are directly related to it based on a citation graph. This procedure makes the LLM input for idea generation more manageable and coherent. In addition, we operationalize the selection process of the core paper and its relevant citations with two design choices: 1) the core paper is selected based on its citation count (e.g., exceeding 100 over 3 months) typically indicating high impact; 2) its relevant papers (which may be potentially numerous) are further narrow-downed based on their similarities of abstracts with the core paper, ensuring a more focused and relevant set of related work.

However, despite the simplicity and intuitiveness of this idea generation approach, there exists one major limitation. This approach relies exclusively on a set of given papers (the core paper and its citations); however, since scientific knowledge is not confined to specific studies but rather accumulates across a wide range of publications (across various fields), we should ideally harness this extensive, interconnected, and relevant scientific knowledge in our method for research idea generation.

Entity-Centric Knowledge Augmentation To achieve this goal, the next question to answer is how is the knowledge in scientific literature \mathcal{L} extracted, stored, and used effectively. In this work, we view entities as the atomic units of knowledge, which allows for ease of its accumulation over numerous papers in a unified manner across different

disciplines. For example, we can easily extract the term database whenever it appears in any paper, using existing off-the-shelf entity linking methods¹ and then aggregate their linked occurrences into a knowledge store. Then, if the term database is prevalent within the realm of medical science but less so in hematology (which is a subdomain of medical science), the constructed knowledge store captures the relevance between those two domains based on overlapping entities (other than the database) and then offers the term database when formulating the ideas about hematology. In other words, this approach enables providing novel and interdisciplinary insights by leveraging the interconnectedness of entities across various fields.

Formally, we design the knowledge store as a two-dimensional matrix $\mathcal{K} \in \mathcal{R}^{m \times m}$ where m is the total number of unique entities identified and \mathcal{K} is implemented in a sparse format. This knowledge store is constructed by extracting entities over all the available scientific articles in literature \mathcal{L}^2 , which not only counts the co-occurrences between entity pairs within individual papers but also quantifies the count for each entity. In addition, to operationalize entity extraction, we use an existing entity linker EL (Wu et al., 2020) that tags and canonicalizes entities in a specific paper l from \mathcal{L} , formalized as follows: $\mathcal{E}_l = \text{EL}(l)$ where \mathcal{E}_l denotes a multiset of entities (allowing for repetitions) appearing in l^3 . Upon extracting entities \mathcal{E} , to store them into the knowledge store \mathcal{K} , we consider all possible pairs of \mathcal{E} represented as follows: $\{e_i, e_j\}_{(i,j) \in \mathcal{C}(|\mathcal{E}|, 2)}$ where $e \in \mathcal{E}$, which is then recorded into \mathcal{K} .

Given this knowledge store \mathcal{K} , our next goal is to enhance the previous vanilla research idea generation process based on a group of interconnected papers, denoted as follows: $\mathbf{o} = \text{LLM}(\mathcal{T}(\{l_0, l_1, \dots, l_n\}))$. We do this by augmenting the LLM with the relevant entities from \mathcal{K} , which can expand the contextual knowledge – what LLMs can consume – by offering additional knowledge. In other words, this knowledge is not seen in the current group of papers but is relevant to it, identified based on entity (co-)occurrence information stored in \mathcal{K} . Formally, let us define entities extracted from the group of interconnected papers,

¹Entity linking is a process that identifies distinct entities in a text and maps them to entities in a knowledge base.

²As extracting entities on all the articles available is not feasible, we target papers appearing after May 01, 2023.

³Due to the extensive length of scientific publications, the target of our entity extraction is titles and abstracts.

as follows: $\mathcal{E}_{\{l_0, \dots, l_n\}} = \bigcup_{i=0}^n \text{EL}(l_i)$. Then, the probabilistic form of retrieving the top- k relevant external entities can be represented as follows:

$$\text{Ret}(\{l_0, \dots, l_n\}; \mathcal{K}) = \arg \max_{I \subset [m]: |I|=k} \prod P(e_i | \mathcal{E}_{\{l_0, \dots, l_n\}}), \quad (1)$$

where $[m] = \{1, \dots, m\}$ and $e_i \notin \mathcal{E}_{\{l_0, \dots, l_n\}}$. Also, for simplicity, by applying Bayes' rule and assuming that entities are independent, the retrieval operation (Equation 1) can be approximated as follows:

$$\arg \max_{I \subset [m]: |I|=k} \prod_{e_j \in \mathcal{E}_{\{l_0, \dots, l_n\}}} P(e_j | e_i) \times P(e_i), \quad (2)$$

where $P(e_j | e_i)$ and $P(e_i)$ can be derived from values in the two-dimensional \mathcal{K} , suitably normalized. Hereafter, the instantiation of research proposal generation augmented with relevant entity-centric knowledge is represented as follows: $\mathbf{o} = \text{LLM}(\mathcal{T}(\{l_0, l_1, \dots, l_n\}, \text{Ret}(\{l_0, \dots, l_n\}; \mathcal{K})))$. We call this knowledge-augmented LLM-powered idea generation approach *ResearchAgent*, and provide the templates to instantiate it in Tables 4, 5, and 6.

3.3 Iterative Research Idea Refinements with Human Preference-Aligned LLM Agents

We note that attempting to write a full research idea in one go may not be an effective strategy, which does not align with the human practice where drafts are continually improved based on multiple reviews and feedback. Therefore, we propose an iterative enhancement strategy, where the LLM-powered reviewing agents (called *ReviewingAgents*) provide the review and feedback according to specific criteria to validate the generated research ideas.

Specifically, similar to our approach to instantiate *ResearchAgent* with an LLM (LLM) and template (\mathcal{T}), *ReviewingAgents* are instantiated similarly but with different templates (See Tables 7, 8, and 9). Then, with *ReviewingAgents*, each of the generated research ideas (problem, method, and experiment design) is separately evaluated according to its own specific five criteria⁴, which are provided in labels of Figure 2 and detailed in Table 10. In addition, based on the reviews and feedback from *ReviewingAgents*, the *ResearchAgent* further updates the already generated research ideas.

Despite the proficiency of LLMs in the evaluation of machine-generated texts (Zheng et al., 2023; Fu et al., 2023), their judgments on the research ideas may not be aligned with the judgments of

⁴We select the top five criteria which we consider as the most important, and leave exploring others as future work.

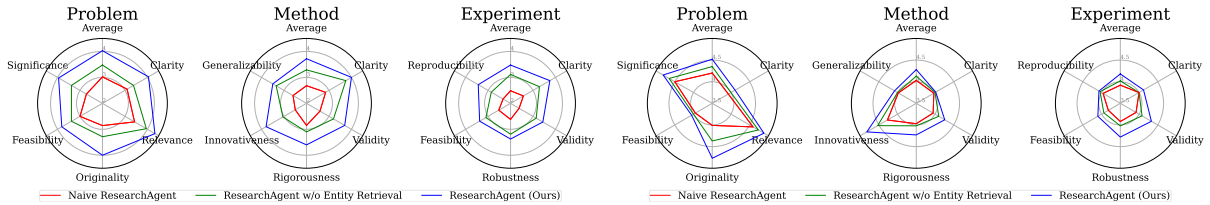


Figure 2: Main results on our research idea generation task with human- (left) and model-based (right) evaluations, where we report the score of each idea (problem, method, or experiment design) based on its own five criteria and their average score.

humans. On the other hand, there are no ground truth reference judgments available, and collecting them to align LLM capabilities is expensive and often infeasible. Ideally, the judgments made by LLMs should be similar to the ones by humans, and we aim to ensure this by automatically generating human preference-aligned evaluation criteria (used for automatic evaluations) with a few human annotations. Specifically, to obtain these human-aligned evaluation criteria, we first collect 10 pairs of the research idea and its score (on a 5-point Likert scale annotated by human researchers having at least 3 papers) on every evaluation criterion. After that, we prompt the LLM with those human-annotated pairs, to induce the detailed descriptions for evaluation criteria that reflect the human preferences, which are then used in ReviewingAgents by including them in evaluation prompt template \mathcal{T} .

4 Experimental Setups

In this section, we describe the datasets, models, evaluation setup, and implementation details.

4.1 Data

The main source to generate research ideas is scientific literature \mathcal{L} , which we obtain from Semantic Scholar Academic Graph API⁵. From this, we select papers appearing after May 01, 2024, because LLMs that we use in our experiments are trained on data from the open web available before this point. Then, we select high-impact papers (that have more than 20 citations) as core papers, mirroring the human researchers’ tendency to leverage influential work, to ensure the high quality of the generated ideas. The resulting data is still very large; therefore, we further randomly sample a subset of 300 papers as core papers (to obtain a reasonably sized benchmark dataset), which means we subsequently generate and evaluate 300 research ideas for each model. The average number of reference papers for each core paper is 87; the abstract of each paper has 2.17 entities on average. The distribution of disciplines for all papers is provided in Figure 7.

⁵<https://www.semanticscholar.org/product/api>

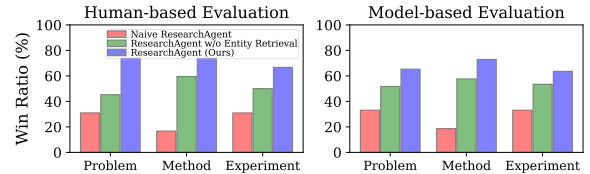


Figure 3: Results of pairwise comparisons between ideas from two of any different approaches, where we report the win ratio.

4.2 Baselines and Our Model

In this work, as we target the novel task of research idea generation, there are no baselines available for direct comparison. Thus, we compare our full ResearchAgent model, which utilizes both references and entities, against ablated variants as follows: 1. **Naive ResearchAgent** – which uses only a core paper to generate research ideas. 2. **ResearchAgent w/o Entity Retrieval** – which uses the core paper and its relevant references without considering entities. 3. **ResearchAgent** – which is our full model that uses the relevant references and entities along with the core paper, to augment LLMs.

4.3 Evaluation Setups

Given that research idea generation is a new task, there are no ground-truth answers to measure the quality of generation. In addition, constructing new pairs of core papers and research ideas is sub-optimal, since there may exist a large number of valid research ideas for each core paper, and this process requires much time, effort and expertise on the part of human researchers. Therefore, we turn to model-based automatic evaluation as well as manual human evaluation to validate different models on our experimental benchmark.

Model-based Evaluation Following the recent trends in using LLMs to judge the quality of output texts (especially in the setting of reference-free evaluations) (Zheng et al., 2023; Fu et al., 2023), we use GPT-4 to judge the quality of research ideas. We note that each of the problem, method, and experiment design is evaluated with five different criteria (See labels of Figure 2 for criteria used for every idea). Then, we ask the evaluation model to either rate the generated idea on a 5-point Likert scale for each criterion or perform pairwise comparisons between two ideas from different models.

Table 1: Results of agreements between two human annotation results and between human and model evaluation results.

Categories	Metrics	Problem	Method	Experiment
Human and Human	Scoring	0.83	0.76	0.67
	Pairwise	0.62	0.62	0.41
Human and Model	Scoring	0.64	0.58	0.49
	Pairwise	0.71	0.62	0.52

We provide the detailed human-induced criteria and prompts used to elicit evaluations in Appendix A.

Human Evaluation Similar to model-based evaluations, we perform human evaluations that involve assigning a score for each criterion and conducting pairwise comparisons between two ideas, with 10 expert annotators. As the generated ideas are knowledge-intensive, it is crucial to select annotators (who are well-versed in the field) and provide them with ideas that are relevant to their field of expertise. Thus, we choose annotators who have authored at least three papers and ask them to judge ideas that are generated from their own papers.

4.4 Implementation Details

We use the GPT-4 (OpenAI, 2023) release from Nov 06 as the basis for all models, which is, notably, reported to be trained with data up to Apr 2023 (meanwhile, the papers used for idea generation appear after May 2023). To extract entities and build the entity-centric knowledge store, we use the off-the-shelf BLINK entity linker (Wu et al., 2020). We provide prompts used to elicit responses for research idea generation in Appendix A.3.

5 Experimental Results and Analyses

We present experimental results and various analyses, showing the effectiveness of ResearchAgent.

Main Results Our main results on scoring with human and model-based evaluations are provided in Figure 2. These demonstrate that our full ResearchAgent outperforms all baselines by large margins on all metrics across the generated problems, methods, and experiment designs (constituting the complete research ideas). Particularly, the full ResearchAgent augmented with relevant entities exhibits strong gains on metrics related to creativity (such as Originality for problems and Innovativeness for methods) since entities may offer novel concepts and views that may not be observable in the group of papers (core paper and its references) used for generating ideas. In addition, the results of pairwise comparisons between two of any models with human and model-based evaluations are reported in Figure 3, on which the full ResearchAgent shows the highest win ratio over its baselines.

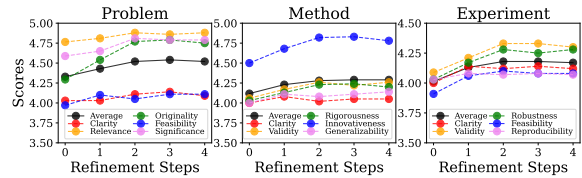


Figure 4: Results with varying the number of refinement steps.

Analysis on Inter-Annotator Agreements To validate the quality and reliability of human annotations, we measure the inter-annotator agreements, where 20% of the generated ideas are evaluated by two humans, and report the results in Table 1. Specifically, for the scoring, we first rank scores from each annotator and measure Spearman’s correlation coefficient (Pirie, 2006) between the ranked scores of two annotators. For the pairwise comparison between two judges, we measure Cohen’s kappa coefficient (Cohen, 1960). As shown in Table 1, we observe that inter-annotator agreement is high, confirming the reliability of our assessments about the quality of generated research ideas.

Analysis on Human-Model Agreements Similar to what we did for the aforementioned inter-annotator agreements, we measure agreements between human-based and model-based evaluations, to ensure the reliability of model-based evaluations. As shown in Table 1, we further confirm that agreements between humans and models are high, indicating that model-based evaluations are a reasonable alternative to judge research idea generation.

Analysis of Refinement Steps To see the effectiveness of iterative refinements of research ideas with ReviewingAgents, in Figure 4, we report the averaged scores on the generated ideas as a function of refinement steps. Based on this, we observe initial improvements in the quality of generated ideas as the number of refinement steps increases. However, the performance becomes saturated after three iterations, which may indicate diminishing returns for subsequent iteration steps.

Ablation on Knowledge Sources Recall that the proposed full ResearchAgent is augmented with two different knowledge sources, namely relevant references and entities. To see the individual contribution of each, we perform an ablation study by either excluding one of the knowledge sources or replacing it with random elements. As shown in Table 2, we observe that each knowledge source contributes to performance improvement. In addition, the performances drop substantially without relevant references, which confirms their importance in generating high-quality research ideas.

Table 2: Results of ablation study on references and entities.

Methods	Problem	Method	Experiment
ResearchAgent	4.52	4.28	4.18
- w/o Entities	4.35	4.13	4.02
- w/ Random Entities	4.41	4.19	4.13
- w/o References	4.26	4.08	3.97
- w/ Random References	4.35	4.16	4.02
- w/o Entities & References	4.20	4.03	3.92

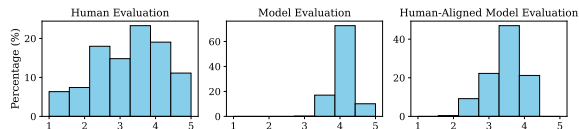


Figure 5: Distributions of model-based evaluation results with and without the human-induced score criteria alignment (middle and right), as well as human evaluation results (left).

Analysis on Human Alignment for Evaluation

Recall that to align judgments from model-based evaluations with actual human preferences, we generated the evaluation criteria based on human evaluation results and used them as the criteria for model-based evaluations. Figure 5 demonstrates the efficacy of this strategy, presenting the score distribution of human evaluation compared with the distributions of model-based evaluations with and without human alignment. We observe that the score distribution of model-based evaluations without human alignment is skewed and far different from the score distribution of human judgments. Yet, after aligning the model-based evaluations with human-induced score criteria, the calibrated distribution more closely resembles the distribution of humans.

Correlation on Citation Counts We further investigate whether a high-impact paper (when used as a core paper) leads to high-quality research ideas. To measure this, we bucketize all papers into three groups by the number of their citations (using it as a proxy for impact), and visualize the average score of each bucket (with model-based evaluations) in Figure 6. We observe that the research ideas generated from high-impact papers are generally of high quality. Additionally, based on the paper distribution (See Figure 7) and for the ease of manual quality check, evaluation criteria for model-based evaluations are induced mainly with computer science papers. To see whether those criteria are applicable to diverse fields, we also compare a correlation between scores of computer science papers and all papers in Figure 6. From this, we observe that the scores increase when the citation increases for both domains, which may support the generalizability of human-preference-induced evaluation criteria.

Analysis using Different LLMs To see how the performance of ResearchAgent changes if an LLM

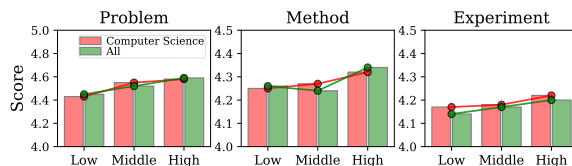


Figure 6: Results with bucketing papers based on citations.

Table 3: Results with different LLMs: GPT-3.5 and GPT-4.0.

LLMs	Models	Problem	Method	Experiment
GPT-4.0	Naive ResearchAgent	4.20	4.03	3.92
	ResearchAgent (Ours)	4.52	4.28	4.18
GPT-3.5	Naive ResearchAgent	3.56	3.56	3.63
	ResearchAgent (Ours)	3.58	3.58	3.60

other than the (most powerful) GPT-4 is used, we conduct an auxiliary analysis instantiating the ResearchAgent with GPT-3.5 (which performs very similarly with the leading open-source LLMs (Touvron et al., 2023)) and present the model-based evaluation results in Table 3. From this, we observe that the performance of ResearchAgent with less capable GPT-3.5 drops significantly, further justifying our choice to not consider weaker LLMs than GPT-4. In addition, the performance differences between the Naive ResearchAgent without knowledge augmentation and the full ResearchAgent become marginal. These results indicate that GPT-3.5 might simply not be capable of capturing complex concepts and their relationships across different scientific papers. This is unsurprising if taken in the context of the emergent abilities of LLMs for complex reasoning (but not in smaller LMs) – a well-known phenomenon (Wei et al., 2022).

6 Conclusion

In this work, we proposed ResearchAgent - a system that aims to accelerate scientific research by automatically generating research ideas, which involves sequential steps of problem identification, method development, and experiment design. In our system, we enhanced LLMs for effective scientific idea generation by leveraging paper relationships over the citation graph and relevant entities extracted and aggregated from numerous papers. Further, we proposed to iteratively refine the generated ideas based on reviews and feedback from LLM-powered multiple reviewing agents, whose evaluation criteria are aligned with human preferences. Through human and model-based evaluations, we showed that ResearchAgent generates research ideas that are more creative, valid, and clear than ones from baselines. We envision ResearchAgent as a collaborative partner (beyond a tool) that strengthens the synergy between researchers and AI in discovering exciting research opportunities.

708 Limitations

709 In this work, we aim to accelerate the first phase
710 of scientific research, demonstrating that the pro-
711 posed ResearchAgent generates useful research
712 ideas. However, there are some areas that future
713 work may improve upon. First of all, recall that
714 we built the entity-centric knowledge store to of-
715 fer fruitful entities during idea generation, and this
716 store is constructed by extracting entities on the
717 titles and abstracts of the limited number of pub-
718 lications (due to the costs of processing them). In
719 addition, the number of entities that we obtain from
720 the BLINK entity linker (Wu et al., 2020) per pa-
721 per may be considered minimal (which is around
722 3). We argue that to build a more comprehensive
723 entity-centric knowledge store, future work may
724 not only extend the content (including the main
725 texts of the publications) and the volume of papers
726 for entity extraction, but also improve the capa-
727 bility of the entity linker itself to more accurately
728 extract scientific terms within the literature. In ad-
729 dition, looking ahead, to truly accelerate the entire
730 scientific research process, experimental validation
731 of the generated research ideas is required, which
732 is a process that is currently time-consuming and
733 demands substantial human efforts. We leave the
734 exploration of this subsequent phase as future work.

735 Ethics Statement

736 We are aware that the ResearchAgent may have the
737 potential to be misused for harmful purposes, such
738 as generating research ideas about new explosives,
739 malicious software, and invasive surveillance tools.
740 Notably, this vulnerability is not unique to our ap-
741 proach but a common challenge faced by existing
742 LLMs that possess significant creative and reason-
743 ing capabilities, occasionally generating content
744 that may be deemed undesirable. Consequently, it
745 underscores the necessity to enhance the robustness
746 and safety of LLMs more broadly.

747 References

748 Microsoft Research AI4Science and Microsoft Azure
749 Quantum. 2023. [The impact of large language mod-
750 els on scientific discovery: a preliminary study using
751 gpt-4.](#) *arXiv preprint arXiv:2311.07361*.

752 Rohan Anil, Sebastian Borgeaud, Yonghui Wu, Jean-
753 Baptiste Alayrac, Jiahui Yu, Radu Soricut, Johan
754 Schalkwyk, Andrew M. Dai, Anja Hauth, Katie Mil-
755 lican, David Silver, Slav Petrov, Melvin Johnson,
756 Ioannis Antonoglou, Julian Schrittwieser, Amelia

Glaese, Jilin Chen, Emily Pitler, Timothy P. Lill-
757 crap, Angeliki Lazaridou, Orhan Firat, James Molloy,
758 Michael Isard, Paul Ronald Barham, Tom Hennig-
759 gan, Benjamin Lee, Fabio Viola, Malcolm Reynolds,
760 Yuanzhong Xu, Ryan Doherty, Eli Collins, Clemens
761 Meyer, Eliza Rutherford, Erica Moreira, Kareem
762 Ayoub, Megha Goel, George Tucker, Enrique Pi-
763 queras, Maxim Krikun, Iain Barr, Nikolay Savinov,
764 Ivo Danihelka, Becca Roelofs, Anaïs White, Anders
765 Andreassen, Tamara von Glehn, Lakshman Yagati,
766 Mehran Kazemi, Lucas Gonzalez, Misha Khalman,
767 Jakub Sygnowski, and et al. 2023. [Gemini: A family
768 of highly capable multimodal models.](#) *arXiv preprint
769 arXiv:2312.11805*. 770

Jinheon Baek, Alham Fikri Aji, and Amir Saffari. 2023. 771
[Knowledge-augmented language model prompting
772 for zero-shot knowledge graph question answering.](#)
773 In *Proceedings of the 1st Workshop on Natural
774 Language Reasoning and Structured Explanations
775 (NLRSE)*, pages 78–106, Toronto, Canada. Associa-
776 tion for Computational Linguistics. 777

Jinheon Baek, Nirupama Chandrasekaran, Silviu
778 Cucerzan, Allen Herring, and Sujay Kumar Jauhar.
779 2024. [Knowledge-augmented large language models
780 for personalized contextual query suggestion.](#) WWW. 781

Andrés M Bran, Sam Cox, Oliver Schilter, Carlo Baldas-
782 sari, Andrew D. White, and Philippe Schwaller. 2023. 783
[Chemcrow: Augmenting large-language models with
784 chemistry tools.](#) 785

Jacob Cohen. 1960. [A coefficient of agreement for
786 nominal scales.](#) *Educational and Psychological Mea-
787 surement*, 20:37 – 46. 788

Michael Fire and Carlos Guestrin. 2019. [Over-
789 optimization of academic publishing metrics: Ob-
790 serving goodhart’s law in action.](#) *GigaScience*, 8. 791

Jinlan Fu, See-Kiong Ng, Zhengbao Jiang, and Pengfei
792 Liu. 2023. [Gptscore: Evaluate as you desire.](#) *arXiv
793 preprint arXiv:2302.04166*. 794

Deep Ganguli, Amanda Askell, Nicholas Schiefer,
795 Thomas I. Liao, Kamile Lukosiute, Anna Chen,
796 Anna Goldie, Azalia Mirhoseini, Catherine Olsson,
797 Danny Hernandez, Dawn Drain, Dustin Li, Eli Tran-
798 Johnson, Ethan Perez, Jackson Kernion, Jamie Kerr,
799 Jared Mueller, Joshua Landau, Kamal Ndousse, Ka-
800 rina Nguyen, Liane Lovitt, Michael Sellitto, Nelson
801 Elhage, Noemí Mercado, Nova DasSarma, Oliver
802 Rausch, Robert Lasenby, Robin Larson, Sam Ringer,
803 Sandipan Kundu, Saurav Kadavath, Scott Johnston,
804 Shauna Kravec, Sheer El Showk, Tamera Lanham,
805 Timothy Telleen-Lawton, Tom Henighan, Tristan
806 Hume, Yuntao Bai, Zac Hatfield-Dodds, Ben Mann,
807 Dario Amodei, Nicholas Joseph, Sam McCandlish,
808 Tom Brown, Christopher Olah, Jack Clark, Samuel R.
809 Bowman, and Jared Kaplan. 2023. [The capacity
810 for moral self-correction in large language models.](#)
811 *arXiv preprint arXiv:2302.07459*. 812

813	Sam Henry and Bridget T. McInnes. 2017. Literature based discovery: Models, methods, and trends . <i>Journal of biomedical informatics</i> , 74:20–32.	868
814		869
815		870
816	Tom Hope, Doug Downey, Daniel S. Weld, Oren Etzioni, and Eric Horvitz. 2023. A computational inflection for scientific discovery . <i>Commun. ACM</i> , 66(8):62–73.	871
817		872
818		873
819		874
820	Qian Huang, Jian Vora, Percy Liang, and Jure Leskovec. 2023. Benchmarking large language models as AI research agents . <i>arXiv preprint arXiv:2310.03302</i> .	875
821		876
822		877
823	Angeliki Lazaridou, Elena Gribovskaya, Wojciech Stokowiec, and Nikolai Grigorev. 2022. Internet-augmented language models through few-shot prompting for open-domain question answering . <i>arXiv preprint arXiv:2203.05115</i> .	878
824		879
825		880
826		881
827		882
828	Nelson F. Liu, Kevin Lin, John Hewitt, Ashwin Paranjape, Michele Bevilacqua, Fabio Petroni, and Percy Liang. 2023. Lost in the middle: How language models use long contexts . <i>Transactions of the Association for Computational Linguistics</i> , 12:157–173.	883
829		884
830		885
831		886
832		887
833	Aman Madaan, Niket Tandon, Prakhar Gupta, Skyler Hallinan, Luyu Gao, Sarah Wiegreffe, Uri Alon, Nouha Dziri, Shrimai Prabhumoye, Yiming Yang, Shashank Gupta, Bodhisattwa Prasad Majumder, Katherine Hermann, Sean Welleck, Amir Yazdanbakhsh, and Peter Clark. 2023. Self-refine: Iterative refinement with self-feedback . In <i>Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023</i> .	888
834		889
835		890
836		891
837		892
838		893
839		894
840		895
841		896
842		897
843		898
844	R.K. Nadkarni, David Wadden, Iz Beltagy, Noah A. Smith, Hannaneh Hajishirzi, and Tom Hope. 2021. Scientific language models for biomedical knowledge base completion: An empirical study . <i>ArXiv</i> , abs/2106.09700.	899
845		900
846		901
847		902
848		903
849	OpenAI. 2023. GPT-4 technical report . <i>arXiv preprint arXiv:2303.08774</i> .	904
850		905
851		906
852		907
853	W. Pirie. 2006. <i>Spearman Rank Correlation Coefficient</i> , volume 8.	908
854		909
855		910
856		911
857	Biqing Qi, Kaiyan Zhang, Haoxiang Li, Kai Tian, Si-hang Zeng, Zhang-Ren Chen, and Bowen Zhou. 2023. Large language models are zero shot hypothesis proposers . <i>arXiv preprint arXiv:2311.05965</i> .	912
858		913
859		914
860		915
861		916
862	Ori Ram, Yoav Levine, Itay Dalmedigos, Dor Muhlgay, Amnon Shashua, Kevin Leyton-Brown, and Yoav Shoham. 2023. In-context retrieval-augmented language models . <i>Transactions of the Association for Computational Linguistics</i> , 11:1316–1331.	917
863		918
864		919
865		920
866		921
867		922
		923
		924
		925
		926
		927
		928
		929
		930
		931
		932
		933
		934
		935
		936
		937
		938
		939
		940
		941
		942
		943
		944
		945
		946
		947
		948
		949
		950
		951
		952
		953
		954
		955
		956
		957
		958
		959
		960
		961
		962
		963
		964
		965
		966
		967
		968
		969
		970
		971
		972
		973
		974
		975
		976
		977
		978
		979
		980
		981
		982
		983
		984
		985
		986
		987
		988
		989
		990
		991
		992
		993
		994
		995
		996
		997
		998
		999
		1000

925 Shengchao Liu, Peter Van Katwyk, Andreea Deac,
926 Anima Anandkumar, Karianne Bergen, Carla P.
927 Gomes, Shirley Ho, Pushmeet Kohli, Joan Lasenby,
928 Jure Leskovec, Tie-Yan Liu, Arjun Manrai, Deb-
929 ora S. Marks, Bharath Ramsundar, Le Song, Jimeng
930 Sun, Jian Tang, Petar Velickovic, Max Welling, Lin-
931 feng Zhang, Connor W. Coley, Yoshua Bengio, and
932 Marinka Zitnik. 2023a. [Scientific discovery in the
933 age of artificial intelligence](#). *Nat.*, 620(7972):47–60.

934 Qingyun Wang, Doug Downey, Heng Ji, and Tom Hope.
935 2023b. [Learning to generate novel scientific direc-
936 tions with contextualized literature-based discovery](#).
937 *arXiv preprint arXiv:2305.14259*.

938 Jason Wei, Yi Tay, Rishi Bommasani, Colin Raf-
939 fel, Barret Zoph, Sebastian Borgeaud, Dani Yo-
940 gatama, Maarten Bosma, Denny Zhou, Donald Met-
941 zler, Ed Huai hsin Chi, Tatsunori Hashimoto, Oriol
942 Vinyals, Percy Liang, Jeff Dean, and William Fedus.
943 2022. [Emergent abilities of large language models](#).
944 *arXiv preprint arXiv:2206.07682*.

945 Sean Welleck, Ximing Lu, Peter West, Faeze Brah-
946 man, Tianxiao Shen, Daniel Khachabi, and Yejin
947 Choi. 2023. [Generating sequences by learning to
948 self-correct](#). In *The Eleventh International Confer-
949 ence on Learning Representations, ICLR 2023, Ki-
950 gali, Rwanda, May 1-5, 2023*. OpenReview.net.

951 Ledell Wu, Fabio Petroni, Martin Josifoski, Sebastian
952 Riedel, and Luke Zettlemoyer. 2020. [Scalable zero-
953 shot entity linking with dense entity retrieval](#). In
954 *Proceedings of the 2020 Conference on Empirical
955 Methods in Natural Language Processing (EMNLP)*,
956 pages 6397–6407, Online. Association for Computa-
957 tional Linguistics.

958 Yike Wu, Nan Hu, Sheng Bi, Guilin Qi, J. Ren, An-
959 huan Xie, and Wei Song. 2023. [Retrieve-rewrite-
960 answer: A kg-to-text enhanced llms framework for
961 knowledge graph question answering](#). *arXiv preprint
962 arXiv:2309.11206*.

963 Zonglin Yang, Xinya Du, Junxian Li, Jie Zheng,
964 Soujanya Poria, and Erik Cambria. 2023. [Large
965 language models for automated open-domain sci-
966 entific hypotheses discovery](#). *arXiv preprint
967 arXiv:2309.02726*.

968 Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan
969 Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin,
970 Zhuohan Li, Dacheng Li, Eric P. Xing, Haotong
971 Zhang, Joseph Gonzalez, and Ion Stoica. 2023. [Judg-
972 ing llm-as-a-judge with mt-bench and chatbot arena](#).
973 *arXiv preprint arXiv:2306.05685*.

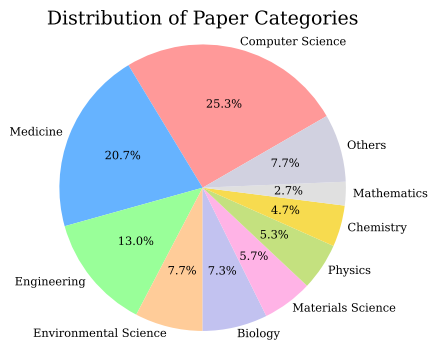


Figure 7: Visualization of the distribution of disciplines for all core papers, selected for research idea generation.

A Additional Experimental Details

In this section, we provide additional details on experiments, including datasets, human evaluation setups, prompts (used for research idea generation and validation), and human-induced criteria.

A.1 Data Statistics

We visualize a distribution of core paper categories used for idea generation in Figure 7, where the categories are obtained from Semantic Scholar API⁶. From this, we find that the top 3 categories are computer science, medicine, and engineering.

A.2 Details on Human Evaluation

To conduct evaluations with human judges, we recruited 10 annotators. They are graduate school students from the United States and South Korea, majoring in computer science, medicine, and biology, each with a minimum of 3 published papers. They were provided with a 6-page guideline document, which includes the task instruction and annotation examples. In addition, they were compensated at a rate of \$22.20 per hour. On average, within an hour, they evaluated 3 sets of research ideas, with each set comprising three sub-ideas (problem, method, and experiment design) from three different approaches (i.e., a total of 9 ideas for one hour). We perform three rounds of human evaluations with refinements in between, and, due to the cost associated with human annotations, we are able to fully evaluate a total of 150 ideas.

A.3 Prompts for Ideas Generation

We provide the prompts used to elicit the idea generations from our full ResearchAgent, specifically for instantiating problem identification, method development, and experiment design in Table 4, Table 5, and Table 6, respectively.

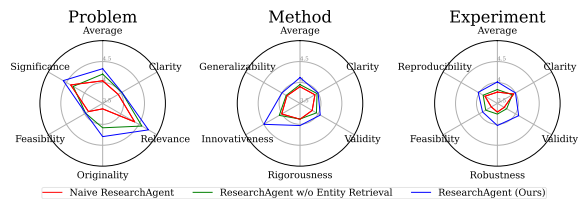


Figure 8: Results on our research idea generation task with model-based evaluation, where we exclude refinement steps.

A.4 Prompts for Idea Validation

We provide the prompts used to elicit the idea validation from our ReviewingAgents as well as the model-based evaluations, specifically for instantiating problem validation, method validation, and experiment design validation in Table 7, Table 8, and Table 9, respectively. In addition, we provide the criteria used, which are induced by human judgments in the next subsection (Appendix A.5).

A.5 Criteria Induced by Human Judgements

Recall that, to align model-based evaluations with human preferences, we induce the criteria (used for automatic evaluations) with actual human judgments. We note that this is done by prompting GPT-4 with 10 pairs of generated ideas and (randomly selected) human judgments. We provide the resulting criteria for validations of problems, methods, and experiment designs in Table 11, Table 12, and Table 13, respectively.

B Additional Experimental Results

We provide additional experimental results, including comparisons without refinements and examples of the generated research ideas.

B.1 Comparisons without Refinements

To see whether the proposed ResearchAgent is consistently effective even without ReviewingAgents, we show the model-based evaluation results without any refinement steps in Figure 8. From this, we clearly observe that the full ResearchAgent outperforms its variants, demonstrating its effectiveness.

B.2 Examples

We provide examples of generated research ideas (including problems, methods, and experiment designs) in Table 14.

⁶<https://www.semanticscholar.org/product/api>

Table 4: The prompt used in the full instantiation of ResearchAgent for problem identification.

Types	Texts
System Message	<p>You are an AI assistant whose primary goal is to identify promising, new, and key scientific problems based on existing scientific literature, in order to aid researchers in discovering novel and significant research opportunities that can advance the field.</p>
	<p>You are going to generate a research problem that should be original, clear, feasible, relevant, and significant to its field. This will be based on the title and abstract of the target paper, those of {len(references)} related papers in the existing literature, and {len(entities)} entities potentially connected to the research area.</p> <p>Understanding of the target paper, related papers, and entities is essential:</p> <ul style="list-style-type: none"> - The target paper is the primary research study you aim to enhance or build upon through future research, serving as the central source and focus for identifying and developing the specific research problem. - The related papers are studies that have cited the target paper, indicating their direct relevance and connection to the primary research topic you are focusing on, and providing additional context and insights that are essential for understanding and expanding upon the target paper. - The entities can include topics, keywords, individuals, events, or any subjects with possible direct or indirect connections to the target paper or the related studies, serving as auxiliary sources of inspiration or information that may be instrumental in formulating the research problem. <p>Your approach should be systematic:</p> <ul style="list-style-type: none"> - Start by thoroughly reading the title and abstract of the target paper to understand its core focus. - Next, proceed to read the titles and abstracts of the related papers to gain a broader perspective and insights relevant to the primary research topic. - Finally, explore the entities to further broaden your perspective, drawing upon a diverse pool of inspiration and information, while keeping in mind that not all may be relevant.
User Message	<p>I am going to provide the target paper, related papers, and entities, as follows:</p> <p>Target paper title: {paper['title']}</p> <p>Target paper abstract: {paper['abstract']}</p> <p>Related paper titles: {relatedPaper['titles']}</p> <p>Related paper abstracts: {relatedPaper['abstracts']}</p> <p>Entities: {Entities}</p> <p>With the provided target paper, related papers, and entities, your objective now is to formulate a research problem that not only builds upon these existing studies but also strives to be original, clear, feasible, relevant, and significant. Before crafting the research problem, revisit the title and abstract of the target paper, to ensure it remains the focal point of your research problem identification process.</p> <p>Target paper title: {paper['title']}</p> <p>Target paper abstract: {paper['abstract']}</p> <p>Then, following your review of the above content, please proceed to generate one research problem with the rationale, in the format of</p> <p>Problem:</p> <p>Rationale:</p>

Table 5: The prompt used in the full instantiation of ResearchAgent for method development.

Types	Texts
System Message	<p>You are an AI assistant whose primary goal is to propose innovative, rigorous, and valid methodologies to solve newly identified scientific problems derived from existing scientific literature, in order to empower researchers to pioneer groundbreaking solutions that catalyze breakthroughs in their fields.</p>
User Message	<p>You are going to propose a scientific method to address a specific research problem. Your method should be clear, innovative, rigorous, valid, and generalizable. This will be based on a deep understanding of the research problem, its rationale, existing studies, and various entities.</p> <p>Understanding of the research problem, existing studies, and entities is essential:</p> <ul style="list-style-type: none"> - The research problem has been formulated based on an in-depth review of existing studies and a potential exploration of relevant entities, which should be the cornerstone of your method development. - The existing studies refer to the target paper that has been pivotal in identifying the problem, as well as the related papers that have been additionally referenced in the problem discovery phase, all serving as foundational material for developing the method. - The entities can include topics, keywords, individuals, events, or any subjects with possible direct or indirect connections to the existing studies, serving as auxiliary sources of inspiration or information that may be instrumental in method development. <p>Your approach should be systematic:</p> <ul style="list-style-type: none"> - Start by thoroughly reading the research problem and its rationale, to understand your primary focus. - Next, proceed to review the titles and abstracts of existing studies, to gain a broader perspective and insights relevant to the primary research topic. - Finally, explore the entities to further broaden your perspective, drawing upon a diverse pool of inspiration and information, while keeping in mind that not all may be relevant. <p>I am going to provide the research problem, existing studies (target paper & related papers), and entities, as follows:</p> <p>Research problem: {researchProblem} Rationale: {researchProblemRationale} Target paper title: {paper['title']} Target paper abstract: {paper['abstract']} Related paper titles: {relatedPaper['titles']} Related paper abstracts: {relatedPaper['abstracts']} Entities: {Entities}</p> <p>With the provided research problem, existing studies, and entities, your objective now is to formulate a method that not only leverages these resources but also strives to be clear, innovative, rigorous, valid, and generalizable. Before crafting the method, revisit the research problem, to ensure it remains the focal point of your method development process.</p> <p>Research problem: {researchProblem} Rationale: {researchProblemRationale}</p> <p>Then, following your review of the above content, please proceed to propose your method with its rationale, in the format of</p> <p>Method: Rationale:</p>

Table 6: The prompt used in the full instantiation of ResearchAgent for experiment design.

Types	Texts
System Message	<p>You are an AI assistant whose primary goal is to design robust, feasible, and impactful experiments based on identified scientific problems and proposed methodologies from existing scientific literature, in order to enable researchers to systematically test hypotheses and validate groundbreaking discoveries that can transform their respective fields.</p>
	<p>You are going to design an experiment, aimed at validating a proposed method to address a specific research problem. Your experiment design should be clear, robust, reproducible, valid, and feasible. This will be based on a deep understanding of the research problem, scientific method, existing studies, and various entities.</p>
	<p>Understanding of the research problem, scientific method, existing studies, and entities is essential:</p> <ul style="list-style-type: none"> - The research problem has been formulated based on an in-depth review of existing studies and a potential exploration of relevant entities. - The scientific method has been proposed to tackle the research problem, which has been informed by insights gained from existing studies and relevant entities. - The existing studies refer to the target paper that has been pivotal in identifying the problem and method, as well as the related papers that have been additionally referenced in the discovery phase of the problem and method, all serving as foundational material for designing the experiment. - The entities can include topics, keywords, individuals, events, or any subjects with possible direct or indirect connections to the existing studies, serving as auxiliary sources of inspiration or information that may be instrumental in your experiment design.
	<p>Your approach should be systematic:</p> <ul style="list-style-type: none"> - Start by thoroughly reading the research problem and its rationale followed by the proposed method and its rationale, to pinpoint your primary focus. - Next, proceed to review the titles and abstracts of existing studies, to gain a broader perspective and insights relevant to the primary research topic. - Finally, explore the entities to further broaden your perspective, drawing upon a diverse pool of inspiration and information, while keeping in mind that not all may be relevant.
User Message	<p>I am going to provide the research problem, scientific method, existing studies (target paper & related papers), and entities, as follows:</p> <p>Research problem: {researchProblem} Rationale: {researchProblemRationale} Scientific method: {scientificMethod} Rationale: {scientificMethodRationale} Target paper title: {paper['title']} Target paper abstract: {paper['abstract']} Related paper titles: {relatedPaper['titles']} Related paper abstracts: {relatedPaper['abstracts']} Entities: {Entities}</p>
	<p>With the provided research problem, scientific method, existing studies, and entities, your objective now is to design an experiment that not only leverages these resources but also strives to be clear, robust, reproducible, valid, and feasible. Before crafting the experiment design, revisit the research problem and proposed method, to ensure they remain at the center of your experiment design process.</p>
	<p>Research problem: {researchProblem} Rationale: {researchProblemRationale} Scientific method: {scientificMethod} Rationale: {scientificMethodRationale}</p>
	<p>Then, following your review of the above content, please proceed to outline your experiment with its rationale, in the format of</p> <p>Experiment: Rationale:</p>

Table 7: The prompt used in the full instantiation of ReviewingAgent for problem validation.

Types	Texts
System Message	<p>You are an AI assistant whose primary goal is to assess the quality and validity of scientific problems across diverse dimensions, in order to aid researchers in refining their problems based on your evaluations and feedback, thereby enhancing the impact and reach of their work.</p>
User Message	<p>You are going to evaluate a research problem for its {metric}, focusing on how well it is defined in a clear, precise, and understandable manner.</p> <p>As part of your evaluation, you can refer to the existing studies that may be related to the problem, which will help in understanding the context of the problem for a more comprehensive assessment.</p> <ul style="list-style-type: none"> - The existing studies refer to the target paper that has been pivotal in identifying the problem, as well as the related papers that have been additionally referenced in the discovery phase of the problem. <p>The existing studies (target paper & related papers) are as follows: Target paper title: {paper['title']} Target paper abstract: {paper['abstract']} Related paper titles: {relatedPaper['titles']} Related paper abstracts: {relatedPaper['abstracts']}</p> <p>Now, proceed with your {metric} evaluation approach that should be systematic:</p> <ul style="list-style-type: none"> - Start by thoroughly reading the research problem and its rationale, keeping in mind the context provided by the existing studies mentioned above. - Next, generate a review and feedback that should be constructive, helpful, and concise, focusing on the {metric} of the problem. - Finally, provide a score on a 5-point Likert scale, with 1 being the lowest, please ensuring a discerning and critical evaluation to avoid a tendency towards uniformly high ratings (4-5) unless fully justified: {criteria} <p>I am going to provide the research problem with its rationale, as follows: Research problem: {researchProblem} Rationale: {researchProblemRationale}</p> <p>After your evaluation of the above content, please provide your review, feedback, and rating, in the format of Review: Feedback: Rating (1-5):</p>

Table 8: The prompt used in the full instantiation of ReviewingAgent for method validation.

Types	Texts
System Message	<p>You are an AI assistant whose primary goal is to assess the quality and soundness of scientific methods across diverse dimensions, in order to aid researchers in refining their methods based on your evaluations and feedback, thereby enhancing the impact and reach of their work.</p>
User Message	<p>You are going to evaluate a scientific method for its {metric} in addressing a research problem, focusing on how well it is described in a clear, precise, and understandable manner that allows for replication and comprehension of the approach.</p> <p>As part of your evaluation, you can refer to the research problem, and existing studies, which will help in understanding the context of the proposed method for a more comprehensive assessment.</p> <ul style="list-style-type: none"> - The research problem has been used as the cornerstone of the method development, formulated based on an in-depth review of existing studies and a potential exploration of relevant entities. - The existing studies refer to the target paper that has been pivotal in identifying the problem and method, as well as the related papers that have been additionally referenced in the discovery phase of the problem and method. <p>The research problem and existing studies (target paper & related papers) are as follows: Research problem: {researchProblem} Rationale: {researchProblemRationale} Target paper title: {paper['title']} Target paper abstract: {paper['abstract']} Related paper titles: {relatedPaper['titles']} Related paper abstracts: {relatedPaper['abstracts']}</p> <p>Now, proceed with your {metric} evaluation approach that should be systematic:</p> <ul style="list-style-type: none"> - Start by thoroughly reading the proposed method and its rationale, keeping in mind the context provided by the research problem, and existing studies mentioned above. - Next, generate a review and feedback that should be constructive, helpful, and concise, focusing on the {metric} of the method. - Finally, provide a score on a 5-point Likert scale, with 1 being the lowest, please ensuring a discerning and critical evaluation to avoid a tendency towards uniformly high ratings (4-5) unless fully justified: {criteria} <p>I am going to provide the proposed method with its rationale, as follows: Scientific method: {scientificMethod} Rationale: {scientificMethodRationale}</p> <p>After your evaluation of the above content, please provide your review, feedback, and rating, in the format of Review: Feedback: Rating (1-5):</p>

Table 9: The prompt used in the full instantiation of ReviewingAgent for experiment design validation.

Types	Texts
System Message	<p>You are an AI assistant whose primary goal is to meticulously evaluate the experimental designs of scientific papers across diverse dimensions, in order to aid researchers in refining their experimental approaches based on your evaluations and feedback, thereby amplifying the quality and impact of their scientific contributions.</p>
User Message	<p>You are going to evaluate an experiment design for its {metric} in validating a scientific method to address a research problem, focusing on how well it is described in a clear, precise, and understandable manner, enabling others to grasp the setup, procedure, and expected outcomes.</p> <p>As part of your evaluation, you can refer to the research problem, scientific method, and existing studies, which will help in understanding the context of the designed experiment for a more comprehensive assessment.</p> <ul style="list-style-type: none"> - The research problem has been formulated based on an in-depth review of existing studies and a potential exploration of relevant entities. - The scientific method has been proposed to tackle the research problem, which has been informed by insights gained from existing studies and relevant entities. - The existing studies refer to the target paper that has been pivotal in identifying the problem, method, and experiment, as well as the related papers that have been additionally referenced in their discovery phases. <p>The research problem, scientific method, and existing studies (target paper & related papers) are as follows: Research problem: {researchProblem} Rationale: {researchProblemRationale} Scientific method: {scientificMethod} Rationale: {scientificMethodRationale} Target paper title: {paper['title']} Target paper abstract: {paper['abstract']} Related paper titles: {relatedPaper['titles']} Related paper abstracts: {relatedPaper['abstracts']}</p> <p>Now, proceed with your {metric} evaluation approach that should be systematic:</p> <ul style="list-style-type: none"> - Start by thoroughly reading the experiment design and its rationale, keeping in mind the context provided by the research problem, scientific method, and existing studies mentioned above. - Next, generate a review and feedback that should be constructive, helpful, and concise, focusing on the {metric} of the experiment. - Finally, provide a score on a 5-point Likert scale, with 1 being the lowest, please ensuring a discerning and critical evaluation to avoid a tendency towards uniformly high ratings (4-5) unless fully justified: {criteria} <p>I am going to provide the designed experiment with its rationale, as follows: Experiment design: {experimentDesign} Rationale: {experimentDesignRationale}</p> <p>After your evaluation of the above content, please provide your review, feedback, and rating, in the format of Review: Feedback: Rating (1-5):</p>

Table 10: The criteria used for evaluating research ideas: problems, methods, and experiment designs.

Types	Criteria	Texts
Problem	Clarity	It assesses whether the problem is defined in a clear, precise, and understandable manner.
	Relevance	It measures whether the problem is pertinent and applicable to the current field or context of study.
	Originality	It evaluates whether the problem presents a novel challenge or unique perspective that has not been extensively explored before.
	Feasibility	It examines whether the problem can realistically be investigated or solved with the available resources and within reasonable constraints.
	Significance	It assesses the importance and potential impact of solving the problem, including its contribution to the field or its broader implications.
Method	Clarity	It assesses whether the method is described in a clear, precise, and understandable manner that allows for replication and comprehension of the approach.
	Validity	It measures the accuracy, relevance, and soundness of the method in addressing the research problem, ensuring that it is appropriate and directly relevant to the objectives of the study.
	Rigorousness	It examines the thoroughness, precision, and consistency of the method, ensuring that the approach is systematic, well-structured, and adheres to high standards of research quality.
	Innovativeness	It evaluates whether the method introduces new techniques, approaches, or perspectives to the research field that differ from standard research practices and advance them in the field.
	Generalizability	It assesses the extent to which the method can be applied to or is relevant for other contexts, populations, or settings beyond the scope of the study.
Experiment	Clarity	It determines whether the experiment design is described in a clear, precise, and understandable manner, enabling others to grasp the setup, procedure, and expected outcomes.
	Validity	It measures the appropriateness and soundness of the experimental design in accurately addressing the research questions or effectively validating the proposed methods, ensuring that the design effectively tests what it is intended to examine.
	Robustness	It evaluates the durability of the experimental design across a wide range of conditions and variables, ensuring that the outcomes are not reliant on a few specific cases and remain consistent across a broad spectrum of scenarios.
	Feasibility	It evaluates whether the experiment design can realistically be implemented with the available resources, time, and technological or methodological constraints, ensuring that the experiment is practical and achievable.
	Reproducibility	It examines whether the information provided is sufficient and detailed enough for other researchers to reproduce the experiment using the same methodology and conditions, ensuring the reliability of the findings.

Table 11: The criteria induced from human judgments for validating the identified problems, which are used to align model-based evaluations with actual human preferences.

Types	Criteria	Texts
Problem	Clarity	1. The problem is presented in a highly ambiguous manner, lacking clear definition and leaving significant room for interpretation or confusion.
		2. The problem is somewhat defined but suffers from vague terms and insufficient detail, making it challenging to grasp the full scope or objective.
		3. The problem is stated in a straightforward manner, but lacks the depth or specificity needed to fully convey the nuances and boundaries of the research scope.
		4. The problem is clearly articulated with precise terminology and sufficient detail, providing a solid understanding of the scope and objectives with minimal ambiguity.
		5. The problem is exceptionally clear, concise, and specific, with every term and aspect well-defined, leaving no room for misinterpretation and fully encapsulating the research scope and aims.
Relevance	Relevance	1. The problem shows almost no relevance to the current field, failing to connect with the established context or build upon existing work.
		2. The problem has minimal relevance, with only superficial connections to the field and a lack of meaningful integration with prior studies.
		3. The problem is somewhat relevant, making a moderate attempt to align with the field but lacking significant innovation or depth.
		4. The problem is relevant and well-connected to the field, demonstrating a good understanding of existing work and offering promising contributions.
		5. The problem is highly relevant, deeply integrated with the current context, and represents a significant advancement in the field.
Originality	Originality	1. The problem exhibits no discernible originality, closely mirroring existing studies without introducing any novel perspectives or challenges.
		2. The problem shows minimal originality, with slight variations from known studies, lacking significant new insights or innovative approaches.
		3. The problem demonstrates moderate originality, offering some new insights or angles, but these are not sufficiently groundbreaking or distinct from existing work.
		4. The problem is notably original, presenting a unique challenge or perspective that is well-differentiated from existing studies, contributing valuable new understanding to the field.
		5. The problem is highly original, introducing a pioneering challenge or perspective that has not been explored before, setting a new direction for future research.
Feasibility	Feasibility	1. The problem is fundamentally infeasible due to insurmountable resource constraints, lack of foundational research, or critical methodological flaws.
		2. The problem faces significant feasibility challenges related to resource availability, existing knowledge gaps, or technical limitations, making progress unlikely.
		3. The problem is feasible to some extent but faces notable obstacles in resources, existing research support, or technical implementation, which could hinder significant advancements.
		4. The problem is mostly feasible with manageable challenges in resources, supported by adequate existing research, and has a clear, achievable methodology, though minor issues may persist.
		5. The problem is highly feasible with minimal barriers, well-supported by existing research, ample resources, and a robust, clear methodology, promising significant advancements.
Significance	Significance	1. The problem shows minimal to no significance, lacking relevance or potential impact in advancing the field or contributing to practical applications.
		2. The problem has limited significance, with a narrow scope of impact and minor contributions to the field, offering little to no practical implications.
		3. The problem demonstrates average significance, with some contributions to the field and potential practical implications, but lacks innovation or broader impact.
		4. The problem is significant, offering notable contributions to the field and valuable practical implications, with evidence of potential for broader impact and advancement.
		5. The problem presents exceptional significance, with groundbreaking contributions to the field, broad and transformative potential impacts, and substantial practical applications across diverse domains.

Table 12: The criteria induced from human judgments for validating the developed methods, which used to align model-based evaluations with actual human preferences.

Types	Criteria	Texts
Method	Clarity	<ol style="list-style-type: none"> 1. The method is explained in an extremely vague or ambiguous manner, making it impossible to understand or replicate the approach without additional information or clarification. 2. The method is described with some detail, but significant gaps in explanation or logic leave the reader with considerable confusion and uncertainty about how to apply or replicate the approach. 3. The method is described with sufficient detail to understand the basic approach, but lacks the precision or specificity needed to fully replicate or grasp the nuances of the methodology without further guidance. 4. The method is clearly and precisely described, with most details provided to allow for replication and comprehension, though minor areas may benefit from further clarification or elaboration. 5. The method is articulated in an exceptionally clear, precise, and detailed manner, enabling straightforward replication and thorough understanding of the approach with no ambiguities.
		<ol style="list-style-type: none"> 1. The method shows a fundamental misunderstanding of the research problem and lacks any credible alignment with established scientific principles or relevant studies. 2. The method partially addresses the research problem but exhibits significant flaws in its scientific underpinning, making its validity questionable despite some alignment with existing literature. 3. The method adequately addresses the research problem but with some limitations in its scientific validity, showing a mix of strengths and weaknesses in its alignment with related studies. 4. The method effectively addresses the research problem, demonstrating a strong scientific basis and sound alignment with existing literature, albeit with minor areas for improvement. 5. The method exemplifies an exceptional understanding of the research problem, grounded in a robust scientific foundation, and shows exemplary integration and advancement of existing studies' findings.
	Rigorousness	<ol style="list-style-type: none"> 1. The method demonstrates a fundamental lack of systematic approach, with significant inconsistencies and inaccuracies in addressing the research problem, showing a disregard for established research standards. 2. The method shows a minimal level of systematic effort but is marred by notable inaccuracies, lack of precision, and inconsistencies that undermine the rigorousness of the method in tackling the research problem. 3. The method exhibits an average level of systematic structure and adherence to research standards but lacks the thoroughness, precision, and consistency required for a rigorous scientific inquiry. 4. The method is well-structured and systematic, with a good level of precision and consistency, indicating a strong adherence to research standards, though it falls short of exemplifying the highest level of rigorousness. 5. The method exemplifies exceptional rigorousness, with outstanding thoroughness, precision, and consistency in its systematic approach, setting a benchmark for high standards in scientific research quality.
		<ol style="list-style-type: none"> 1. The method introduces no novel elements, fully relying on existing techniques without any attempt to modify or adapt them for the specific research problem, showing a lack of innovativeness. 2. The method shows minimal innovation, with only slight modifications to existing techniques that do not substantially change or improve the approach to the research problem. 3. The method demonstrates moderate innovativeness, incorporating known techniques with some new elements or combinations that offer a somewhat fresh approach to the research problem but fall short of a significant breakthrough. 4. The method is highly innovative, introducing new techniques or novel combinations of existing methods that significantly differ from standard practices, offering a new perspective or solution to the research problem. 5. The method represents a groundbreaking innovation, fundamentally transforming the approach to the research problem with novel techniques or methodologies that redefine the field's standard practices.
	Innovativeness	<ol style="list-style-type: none"> 1. The method shows no adaptability, failing to extend its applicability beyond its original context or dataset, showing a complete lack of generalizability. 2. The method demonstrates minimal adaptability, with limited evidence of potential applicability to contexts slightly different from the original. 3. The method exhibits some level of adaptability, suggesting it could be applicable to related contexts or datasets with modifications. 4. The method is adaptable and shows evidence of applicability to a variety of contexts or datasets beyond the original. 5. The method is highly adaptable, demonstrating clear evidence of broad applicability across diverse contexts, populations, and settings.
		Generalizability

Table 13: The criteria induced from human judgments for validating the experiment designs, which are used to align model-based evaluations with actual human preferences.

Types	Criteria	Texts
	Clarity	<ol style="list-style-type: none"> 1. The experiment design is extremely unclear, with critical details missing or ambiguous, making it nearly impossible for others to understand the setup, procedure, or expected outcomes. 2. The experiment design lacks significant clarity, with many important aspects poorly explained or omitted, challenging others to grasp the essential elements of the setup, procedure, or expected outcomes. 3. The experiment design is moderately clear, but some aspects are not detailed enough, leaving room for interpretation or confusion about the setup, procedure, or expected outcomes. 4. The experiment design is mostly clear, with most aspects well-described, allowing others to understand the setup, procedure, and expected outcomes with minimal ambiguity. 5. The experiment design is exceptionally clear, precise, and detailed, enabling easy understanding of the setup, procedure, and expected outcomes, with no ambiguity or need for further clarification.
	Validity	<ol style="list-style-type: none"> 1. The experiment design demonstrates a fundamental misunderstanding of the research problem, lacks alignment with scientific methods, and shows no evidence of validity in addressing the research questions or testing the proposed methods. 2. The experiment design has significant flaws in its approach to the research problem and scientific method, with minimal or questionable evidence of validity, making it largely ineffective in addressing the research questions or testing the proposed methods. 3. The experiment design is generally aligned with the research problem and scientific method but has some limitations in its validity, offering moderate evidence that it can somewhat effectively address the research questions or test the proposed methods. 4. The experiment design is well-aligned with the research problem and scientific method, providing strong evidence of validity and effectively addressing the research questions and testing the proposed methods, despite minor limitations. 5. The experiment design excellently aligns with the research problem and scientific method, demonstrating robust evidence of validity and outstandingly addressing the research questions and testing the proposed methods without significant limitations.
Experiment	Robustness	<ol style="list-style-type: none"> 1. The experiment design demonstrates a fundamental lack of understanding of the scientific method, with no evidence of durability or adaptability across varying conditions, leading to highly unreliable and non-replicable results. 2. The experiment design shows minimal consideration for robustness, with significant oversights in addressing variability and ensuring consistency across different scenarios, resulting in largely unreliable outcomes. 3. The experiment design adequately addresses some aspects of robustness but lacks comprehensive measures to ensure durability and consistency across a wide range of conditions, leading to moderate reliability. 4. The experiment design incorporates a solid understanding of robustness, with clear efforts to ensure the experiment's durability and consistency across diverse conditions, though minor improvements are still possible for optimal reliability. 5. The experiment design exemplifies an exceptional commitment to robustness, with meticulous attention to durability and adaptability across all possible conditions, ensuring highly reliable and universally applicable results.
	Feasibility	<ol style="list-style-type: none"> 1. The experiment design is fundamentally unfeasible, with insurmountable resource, time, or technological constraints that make implementation virtually impossible within the proposed framework. 2. The experiment design faces significant feasibility challenges, including major resource, time, or technological limitations, that heavily compromise its practical execution and likelihood of success. 3. The experiment design is somewhat feasible, with moderate constraints on resources, time, or technology that could be addressed with adjustments, though these may not guarantee success. 4. The experiment design is largely feasible, with minor resource, time, or technological limitations that can be effectively managed or mitigated, ensuring a high probability of successful implementation. 5. The experiment design is highly feasible, with no significant constraints on resources, time, or technology, indicating that it can be implemented smoothly and successfully within the proposed framework.
	Reproducibility	<ol style="list-style-type: none"> 1. The experiment design lacks critical details, making it virtually impossible for other researchers to replicate the study under the same conditions or methodologies. 2. The experiment provides some essential information but omits significant details needed for replication, leading to considerable ambiguity in methodology or conditions. 3. The experiment design includes sufficient details for replication, but lacks clarity or completeness in certain areas, posing challenges for seamless reproducibility. 4. The experiment is well-documented with clear, detailed instructions and methodologies that allow for consistent replication, albeit with minor areas for improvement. 5. The experiment design is exemplary in its clarity, detail, and comprehensiveness, ensuring that other researchers can precisely and effortlessly replicate the study under identical conditions and methodologies.

Table 14: The examples of research idea generation results from the proposed full ResearchAgent.

Index	Types	Texts
1	Input	<p>Title: Knowledge-Augmented Language Model Prompting for Zero-Shot Knowledge Graph Question Answering</p> <p>Abstract: Large Language Models (LLMs) are capable of performing zero-shot closed-book question answering tasks, based on their internal knowledge stored in parameters during pre-training. However, such internalized knowledge might be insufficient and incorrect, which could lead LLMs to generate factually wrong answers. Furthermore, fine-tuning LLMs to update their knowledge is expensive. To this end, we propose to augment the knowledge directly in the input of LLMs. Specifically, we first retrieve the relevant facts to the input question from the knowledge graph based on semantic similarities between the question and its associated facts. After that, we prepend the retrieved facts to the input question in the form of the prompt, which is then forwarded to LLMs to generate the answer. Our framework, Knowledge-Augmented language model PromptING (KAPING), requires no model training, thus completely zero-shot. We validate the performance of our KAPING framework on the knowledge graph question answering task, that aims to answer the user’s question based on facts over a knowledge graph, on which ours outperforms relevant zero-shot baselines by up to 48% in average, across multiple LLMs of various sizes.</p>
		<p>Entities: ['Natural language', 'Learning management system', 'Question answering', 'Wikipedia', 'Artificial intelligence', 'Named-entity recognition', 'Quality assurance', 'Long short-term memory', 'English language', 'Google Books', 'Natural-language programming', 'LMS color space', 'GitHub', 'Logic learning machine', 'Machine translation', 'DBpedia', 'Integrated library system', 'Spanish language', 'Neural machine translation', 'ArXiv', 'Artificial general intelligence', 'Natural-language understanding', 'Nonlinear programming', 'Codex', 'LoRa', 'Adhesive', 'Common Crawl', 'ROUGE (metric)', 'Transformer', 'Tomas Mikolov']</p>
	Problem	Developing a Multimodal Knowledge-Aware Prompting System for Multilingual Zero-Shot Question Answering across Structured and Unstructured Data Sources

Continued on the next page

Table 14 – Continued from the previous page

Index	Types	Texts
		<p>The further refined and enhanced method for developing a multimodal knowledge-aware prompting system for multilingual zero-shot question answering (MMKP-ZSQA) incorporates advanced refinements in the following components to address the requirements for innovativeness, clarity, and generalizability:</p> <ol style="list-style-type: none"> 1. Interconnected Multimodal Knowledge Orchestration (IMKO): <ul style="list-style-type: none"> - IMKO is an advanced evolution of SMKI, incorporating state-of-the-art language and vision models to handle data verbalization with context-aware disambiguation methods. - Techniques like attention mechanisms are employed to weigh the relevance of different data pieces when integrated, revealing how textual and visual data augment each other in a synthesized representation. 2. Linguistically Inclusive Retrieval Engine (LIRE): <ul style="list-style-type: none"> - LIRE extends EMKA with an emphasis on semantic understanding, using transformer-based models trained on diverse datasets, including idiomatic and cultural nuances across languages. - Specific algorithms to handle linguistic phenomena such as code-switching and transliteration are included, enhancing the application to a broader set of languages and dialects. 3. Prompt Learning and Optimization Nexus (PLON): <ul style="list-style-type: none"> - Building on AMPL, PLON focuses on creating a library of optimized prompts categorized by linguistic features and data modalities, using Bayesian optimization algorithms. - It includes domain adaptation techniques and a wider array of meta-learning strategies with case studies for high and low-resource languages, and outlines theoretical frameworks for their implementation. 4. Cross-Modal Integrative Learning System (C-MILS): <ul style="list-style-type: none"> - C-MILS advances SCAS by detailing the use of multi-head attention across modalities for effective and scalable reasoning, with exemplifications on how each modality can enhance comprehension synergistically. - The component now incorporates a layer of abstraction to distill knowledge into a modality-agnostic format, aiding reasoning and facilitating interpretability across languages and data types. 5. User-Centric Adaptation and Privacy Framework (UCAPF): <ul style="list-style-type: none"> - UCAPF enriches ALCUM with a user-centric design, focusing on interaction protocols that describe user engagement cycles, feedback loops, and privacy-preserving active learning paradigms, with practical workflows and pseudocode. - Aligns with GDPR and other privacy frameworks to fortify trust and address the ethical use of data in dynamic learning environments. 6. Global Evaluation and Reflective Testing System (GERTS): <ul style="list-style-type: none"> - GERTS expands CEB by presenting a multi-tiered validation and reflection methodology to adjust system components based on a diverse set of metrics, including fairness, bias, interpretability, and computational efficiency. - Offers a structured approach to cross-cultural evaluation, including the use of demographic and regional diversity in forming test cohorts.
	Method	
		<p>The experiment, named "Refined Experiment for Multimodal Knowledge-Aware Prompting System for Multilingual Zero-shot Question Answering (RE-MKP-ZSQA)", aims to methodically develop and validate an advanced AI system. The experiment is streamlined to address feasibility, clarity, and reproducibility concerns while upholding robustness and validity by adhering to the following refined phases:</p> <ol style="list-style-type: none"> 1. Detailed System Implementation Plan: <ul style="list-style-type: none"> - Provide a publicly accessible project roadmap with specific milestones, resource allocation, and timelines. 2. Dataset Curation with Clear Guidelines: <ul style="list-style-type: none"> - Publish precise annotation guidelines with strategies to prevent bias. - Document the dataset assembly process, including source selection and data processing procedures. 3. Transparent System Training: <ul style="list-style-type: none"> - Offer a detailed training protocol with hyperparameters, optimization strategies, and Bayesian optimization processes used in PLON. 4. Structured Zero-Shot Evaluation: <ul style="list-style-type: none"> - Outline evaluation metrics derived from GERTS with benchmark datasets to test zero-shot capabilities. 5. Clearer Interdisciplinary Evaluation Protocol: <ul style="list-style-type: none"> - Specify the composition of the evaluation committee, criteria for assessments, and methods for integrating the feedback. 6. Iterative Improvement with Validation Metrics: <ul style="list-style-type: none"> - Describe statistical methods for reflective assessment and continuous improvement, aligned with multi-tiered GERTS methodology. 7. User-Centric Design and Privacy Compliance Evaluation: <ul style="list-style-type: none"> - Structure user studies with targeted data points to assess usability and cultural adaptability. - Outline privacy compliance protocols to adhere to international standards. 8. Detailed Global Scalability Evaluation Method: <ul style="list-style-type: none"> - Define evaluation metrics for scalability tests and describe diverse infrastructural setups. 9. Enhanced Reporting for Reproducibility: <ul style="list-style-type: none"> - Commit to creating a comprehensive report with precise specifications, configurations, and instructions for replication purposes. - Utilize GitHub for version-controlled deposition of code and datasets, and arXiv for openly accessible experiment protocols and findings.
	Experiment	

Continued on the next page

Table 14 – Continued from the previous page

Index	Types	Texts
		<p>Title: Test-Time Self-Adaptive Small Language Models for Question Answering</p> <p>Abstract: Recent instruction-finetuned large language models (LMs) have achieved notable performances in various tasks, such as question-answering (QA). However, despite their ability to memorize a vast amount of general knowledge across diverse tasks, they might be suboptimal on specific tasks due to their limited capacity to transfer and adapt knowledge to target tasks. Moreover, further finetuning LMs with labeled datasets is often infeasible due to their absence, but it is also questionable if we can transfer smaller LMs having limited knowledge only with unlabeled test data. In this work, we show and investigate the capabilities of smaller self-adaptive LMs, only with unlabeled test data. In particular, we first stochastically generate multiple answers, and then ensemble them while filtering out low-quality samples to mitigate noise from inaccurate labels. Our proposed self-adaption strategy demonstrates significant performance improvements on benchmark QA datasets with higher robustness across diverse prompts, enabling LMs to stay stable.</p> <p>Entities: ['Codex', 'Natural language', 'English language', 'United States', 'Question answering', 'Natural-language programming', 'GTRI Information and Communications Laboratory', 'Artificial intelligence', 'LoRa', 'Llama', 'Python (programming language)', 'Learning management system', 'Natural language processing', 'Reinforcement learning', 'LMS color space', 'Wikipedia', 'GitHub', 'Natural-language understanding', 'London, Midland and Scottish Railway', 'Integrated library system', 'Language model', 'Chinese language', 'Lumen (unit)', 'Spanish language', 'English Wikipedia', 'Logic learning machine', 'Gradient descent', 'Alternative public offering', 'Technology transfer', 'Dialogue system']</p>
	Input	
	Problem	<p>Developing a Scalable, Domain-Adaptive Test-Time Training Protocol for Low-Resource Language QA Using Small Language Models</p>
2		<ol style="list-style-type: none"> 1. Selection of Scalable Compact Language Models (CLMs): Identify and evaluate existing CLMs suitable for adaptation, emphasizing models with minimal computational requirements. 2. Creation of a Multilingual Test-Time Training (TTT) Framework: Develop a TTT protocol that enables CLMs to adapt to new domains and languages during the inference phase, leveraging unsupervised learning techniques and pseudo-label generation. 3. Synthetic and Unsupervised Data Generation: Utilize a combination of unsupervised and synthetic data generation methods to produce multilingual QA pairs, employing techniques such as back-translation and context-based question synthesis. 4. Domain-Adaptive Mechanisms: Introduce domain-adaptive components, including feature adaptation layers and meta-learning algorithms, which tailor the model's behavior to new contexts and languages at test time. 5. Incremental Language Addition and Dominance Assessment: Start with a subset of linguistically diverse, low-resource languages. Evaluate domain adaptability for each language via an iterative process, ensuring models learn to prioritize resource efficiency. 6. Model Robustness and Generalization: Perform robustness tuning (RT) to prepare models for unforeseen linguistic variations and conduct thorough evaluations across multiple domains to ensure models can generalize their learning effectively. 7. Human-In-The-Loop Evaluation: Conduct evaluations with native speakers and domain experts to validate the relevance and accuracy of the QA outputs, incorporating feedback into the iterative training process. 8. Open-Sourcing and Community Collaboration: Make the TTT protocol, trained models, and evaluation benchmarks publicly available for the research community, fostering collaboration and further innovation.
	Method	

Continued on the next page

Table 14 – Continued from the previous page

Index	Types	Texts
		<p>1. Selection and Preparation:</p> <ul style="list-style-type: none"> - Identify potential compact language models (CLMs) suitable for domain adaptation and test-time training, focusing on those with minimal computational requirement and the ability to be fine-tuned or adapted in an unsupervised manner. - Prepare a diverse set of low-resource languages and corresponding text corpora, ensuring linguistic diversity and sociocultural significance. Select benchmark datasets for these languages if available. <p>2. Training and Adaptation Procedure:</p> <ul style="list-style-type: none"> - Create a Test-Time Training (TTT) framework that allows selected CLMs to adapt to various domains in the selected low-resource languages during the inference phase. - Implement unsupervised learning techniques and pseudo-label generation to produce QA pairs, utilizing back-translation and context-based question synthesis to generate synthetic datasets for languages with limited or no available QA datasets. - Integrate domain-adaptive components and meta-learning algorithms into the CLMs to enable domain-specific adaptations at test time. <p>3. Iterative Evaluation and Refinement:</p> <ul style="list-style-type: none"> - Begin adaptation and training with a single low-resource language and gradually add additional languages, monitoring the domain adaptability and model performance metrics after each addition. - Perform robustness tuning and cross-domain evaluations for each CLM and language adaptation to ensure generalizability and prevent overfitting. <p>4. Human-In-The-Loop Assessment:</p> <ul style="list-style-type: none"> - Enlist native speakers and domain experts to evaluate the relevance and accuracy of the model's QA outputs for each language. - Incorporate feedback into the iterative training process, refining and re-adapting the models accordingly. <p>5. Open-Sourcing and Community Feedback:</p> <ul style="list-style-type: none"> - Make the TTT protocol, adaptive CLMs, evaluation benchmarks, and any synthetic datasets publicly available for the research community. <p>6. Experiment Monitoring and Documentation:</p> <ul style="list-style-type: none"> - Record all the parameters, datasets, model configurations, and evaluation metrics meticulously to ensure robustness and reproducibility. - Document any challenges faced, unexpected results, or adaptations made during the experiment for open-sourcing purposes. <p>7. Data Analysis and Reporting:</p> <ul style="list-style-type: none"> - Analyze the collected performance data quantitatively, using appropriate statistical methods to compare with non-adaptive baselines. - Report qualitative findings from human-in-the-loop evaluations, interpreting the implications for language model performance in low-resource language domains.
3	Input	<p>Title: Whole-brain annotation and multi-connectome cell typing quantifies circuit stereotypy in <i>Drosophila</i></p> <p>Abstract: The fruit fly <i>Drosophila melanogaster</i> combines surprisingly sophisticated behaviour with a highly tractable nervous system. A large part of the fly's success as a model organism in modern neuroscience stems from the concentration of collaboratively generated molecular genetic and digital resources. As presented in our FlyWire companion paper¹, this now includes the first full brain connectome of an adult animal. Here we report the systematic and hierarchical annotation of this 130,000-neuron connectome including neuronal classes, cell types and developmental units (hemilineages). This enables any researcher to navigate this huge dataset and find systems and neurons of interest, linked to the literature through the Virtual Fly Brain database². Crucially, this resource includes 4,552 cell types. 3,094 are rigorous consensus validations of cell types previously proposed in the "hemibrain" connectome³. In addition, we propose 1,458 new cell types, arising mostly from the fact that the FlyWire connectome spans the whole brain, whereas the hemibrain derives from a subvolume. Comparison of FlyWire and the hemibrain showed that cell type counts and strong connections were largely stable, but connection weights were surprisingly variable within and across animals. Further analysis defined simple heuristics for connectome interpretation: connections stronger than 10 unitary synapses or providing >1% of the input to a target cell are highly conserved. Some cell types showed increased variability across connectomes: the most common cell type in the mushroom body, required for learning and memory, is almost twice as numerous in FlyWire as the hemibrain. We find evidence for functional homeostasis through adjustments of the absolute amount of excitatory input while maintaining the excitation-inhibition ratio. Finally, and surprisingly, about one third of the cell types proposed in the hemibrain connectome could not yet be reliably identified in the FlyWire connectome. We therefore suggest that cell types should be defined to be robust to inter-individual variation, namely as groups of cells that are quantitatively more similar to cells in a different brain than to any other cell in the same brain. Joint analysis of the FlyWire and hemibrain connectomes demonstrates the viability and utility of this new definition. Our work defines a consensus cell type atlas for the fly brain and provides both an intellectual framework and open source toolchain for brain-scale comparative connectomics.</p> <p>Entities: ['Virtual Fly Brain', 'Central nervous system', 'Transposable element', 'SUMO protein', 'Kenyon cell', 'Romani people', 'Induced stem cells', 'Ventral nerve cord', 'FlyBase', 'Parkinson's disease', 'Virtual Network Computing', 'P element', 'Piwi-interacting RNA', 'Drosophila Genetic Reference Panel', 'Bateson–Dobzhansky–Muller model', 'J. B. S. Haldane', 'ATG7', 'Haldane's rule', 'Oxford Nanopore Technologies', 'Drosophila mauritiana', 'Germline', 'PINK1', 'Migratory locust', 'CRISPR', 'Helicobacter', 'GINS (protein complex)', 'Parkin (ligase)', 'Lepidoptera', 'Illumina, Inc.', 'Drosophila']</p>
	Problem	Investigating the Functional Implications of Connectome Variability in <i>Drosophila</i> 's Learning and Memory Circuits Across Different Environmental and Genetic Contexts

Continued on the next page

Table 14 – Continued from the previous page

Index	Types	Texts
		<p>The proposed method involves a multi-tiered approach that integrates connectomics, behavioral assays, genetic manipulation, and computational modeling to investigate the functional implications of connectome variability in <i>Drosophila</i>'s learning and memory circuits. The method consists of the following steps:</p> <ol style="list-style-type: none"> 1. Connectome Mapping and Variability Analysis: <ol style="list-style-type: none"> a. Utilize the Virtual Fly Brain database to identify and compare individual connectomes, focusing on the mushroom body. b. Quantify the variability in connection weights and cell type counts using statistical methods and machine learning algorithms to identify patterns of variability. 2. Behavioral Assays: <ol style="list-style-type: none"> a. Design a series of learning and memory tasks for <i>Drosophila</i>, such as olfactory conditioning or visual pattern recognition. b. Test groups of flies with known connectome profiles under controlled environmental conditions to establish baseline behavioral data. 3. Environmental and Genetic Perturbations: <ol style="list-style-type: none"> a. Expose different groups of flies to varied learning paradigms and sensory inputs to create environmental perturbations. b. Use CRISPR-Cas9 technology to introduce targeted mutations in genes like PINK1 or Parkin, creating genetic perturbations. c. Assess the impact of these perturbations on connectome structure using high-resolution imaging and reconstruction techniques. 4. Transcriptomic and Spatial Analysis: <ol style="list-style-type: none"> a. Apply single-cell RNA sequencing and spatial transcriptomics to profile gene expression changes in response to environmental and genetic perturbations. b. Correlate transcriptomic data with connectome changes to identify molecular pathways associated with structural and functional plasticity. 5. Computational Modeling and Network Analysis: <ol style="list-style-type: none"> a. Develop computational models to simulate the effects of connectome variability on neural circuit function. b. Use network analysis tools to explore information flow and circuit dynamics, integrating data from related papers on olfactory projection neurons and information flow in the olfactory system. 6. Synthesis and Validation: <ol style="list-style-type: none"> a. Integrate findings from behavioral assays, connectome mapping, transcriptomic analysis, and computational modeling. b. Validate the proposed models and hypotheses through iterative experimentation and refinement.
Method		<hr/> <ol style="list-style-type: none"> 1. Pre-experimental Setup: <ol style="list-style-type: none"> a. Establish a <i>Drosophila melanogaster</i> breeding program to ensure a consistent supply of genetically similar flies for the experiments. b. Develop a standardized protocol for rearing flies under controlled conditions to minimize pre-experimental variability. c. Select and validate CRISPR-Cas9 constructs for targeted gene editing of PINK1, Parkin, and other genes of interest. d. Train personnel in the use of the Virtual Fly Brain database and relevant computational tools for connectome analysis. 2. Connectome Mapping and Variability Analysis: <ol style="list-style-type: none"> a. Randomly assign individual flies to either a control group or various treatment groups (environmental and genetic perturbations). b. Utilize high-resolution imaging techniques to map the connectomes of flies from each group, with a focus on the mushroom body. c. Apply statistical and machine learning algorithms to quantify and compare the variability in connection weights and cell type counts across groups. 3. Behavioral Assays: <ol style="list-style-type: none"> a. Design and validate a series of learning and memory tasks, such as olfactory conditioning and visual pattern recognition, ensuring tasks are sensitive to subtle differences in performance. b. Test flies from each group in the behavioral tasks and record performance metrics. c. Analyze behavioral data to establish correlations with connectome profiles. 4. Environmental and Genetic Perturbations: <ol style="list-style-type: none"> a. Expose flies to different learning paradigms and sensory inputs to induce environmental perturbations. b. Perform gene editing using CRISPR-Cas9 to create genetic perturbations in the treatment groups. c. Re-map connectomes post-perturbation to assess structural changes. 5. Transcriptomic and Spatial Analysis: <ol style="list-style-type: none"> a. Collect brain tissue from flies post-behavioral assays and perform single-cell RNA sequencing and spatial transcriptomics. b. Analyze transcriptomic data to identify gene expression changes and correlate these with observed connectome and behavioral variations. 6. Computational Modeling and Network Analysis: <ol style="list-style-type: none"> a. Develop computational models to simulate the impact of observed connectome variability on neural circuit function. b. Use network analysis to integrate behavioral, connectomic, and transcriptomic data, focusing on information flow and circuit dynamics. 7. Synthesis and Validation: <ol style="list-style-type: none"> a. Integrate findings across all experimental components to formulate a cohesive understanding of the functional implications of connectome variability. b. Validate models and refine hypotheses through additional targeted experiments, informed by initial findings.
	Experiment	