

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

Anonymous ACL submission

## Abstract

The pace of scientific research, vital for improving human life, is complex, slow, and needs specialized expertise. Meanwhile, novel, impactful research often stems from both a deep understanding of prior work, and a cross-pollination of ideas across domains and fields. To enhance the productivity of researchers, we propose ResearchAgent, which leverages the encyclopedic knowledge and linguistic reasoning capabilities of Large Language Models (LLMs) to assist them in their work. This system automatically defines novel problems, proposes methods and designs experiments, while iteratively refining them based on the feedback from collaborative LLM-powered reviewing agents. Specifically, starting with a core scientific paper, ResearchAgent is augmented not only with relevant publications by connecting information over an academic graph but also entities retrieved from a knowledge store derived from shared underlying concepts mined across numerous papers. Then, mimicking a scientific approach to improving ideas with peer discussions, we leverage multiple LLM-based 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 judgements 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 both human and model-based evaluation results. Our initial foray into AI-mediated scientific research has important implications for the development of future systems aimed at supporting researchers in their ideation and operationalization of novel work.

## 1 Introduction

Scientific research plays a crucial role in driving innovation, advancing knowledge, solving problems, expanding our understanding of the world,

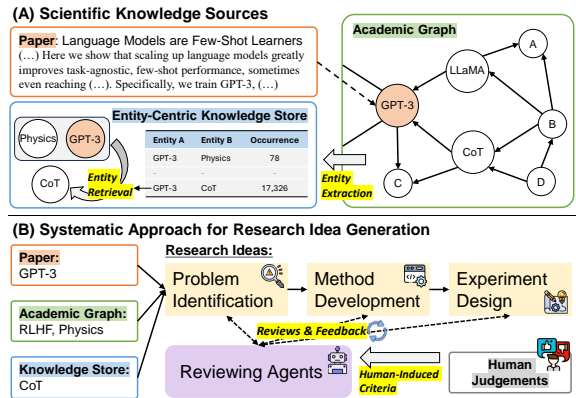


Figure 1: (A) The scientific knowledge used for research idea generation consists of a paper, its relationships over an academic graph, and entities within a knowledge store extracted from numerous papers. (B) Given them, the proposed research idea generation process involves problem identification,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

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 slow, effort-intensive process, which requires reading and synthesizing overwhelming amounts of knowledge over the vast corpus of rapidly growing scientific literature to formulate research ideas, as well as design and perform experimental validations of those ideas. For example, the number of academic papers published per year is more than 7 million (Fire and Guestrin, 2019). Similarly, the process of testing a new pharmaceutical drug requires deep expertise, and is massively expensive and labor-intensive, often taking several years (Vamathevan et al., 2019).

In the meantime, recent Large Language Models (LLMs) (Touvron et al., 2023; OpenAI, 2023; Anil et al., 2023) have shown impressive capabilities in processing and generating text with remark-

067 able accuracy, even outperforming human experts  
068 across diverse specialized domains including math,  
069 physics, history, law, medicine, and ethics. They  
070 are able to process and analyze large volumes of  
071 data at speeds and scales far exceeding human ca-  
072 pabilities, have internalized large swaths of human  
073 knowledge from being trained on virtually the en-  
074 tire web, and can identify patterns, trends, and cor-  
075 relations that may not be immediately apparent to  
076 human researchers (such as the usage of quantum  
077 mechanics in medical imaging or applying psycho-  
078 logical insights in AI). This renders them ideally  
079 poised to become foundational tools to accelerate  
080 the two phases of the scientific research process:  
081 ideation of novel research opportunities, and scien-  
082 tific validation of those research hypotheses.

083 A few recent papers in the domain of LLM-  
084 augmented scientific discovery have focused on the  
085 second phase. Specifically, they attempt (Huang  
086 et al., 2023; AI4Science and Quantum, 2023; Bran  
087 et al., 2023) to mainly accelerate the experimental  
088 validation process, by writing code for machine-  
089 learning models, facilitating the exploration of  
090 chemical spaces, or advancing the simulation of  
091 molecular dynamics. Thus, in this paper, we lever-  
092 age LLMs in the first phase of scientific research  
093 – specifically idea generation, whose key focus is  
094 conceptualizing novel research questions, method-  
095 ologies, and experiments. To the best of our knowl-  
096 edge, our work is the first to leverage and evaluate  
097 the capabilities of LLMs to act as mediators in  
098 scientific idea generation in an open-ended setting.

099 Given our goal to build an LLM-powered Re-  
100 searchAgent, we draw inspiration from how human  
101 researchers position themselves to come up with  
102 novel research ideas. We draw distinctions between  
103 three key components of their workflow: a broad  
104 and deep understanding of related scientific liter-  
105 ature, an encyclopedic view of concepts and how  
106 they relate to one another both within and across  
107 domains, and a community of colleagues on which  
108 to rely for feedback and constructive criticism.

109 We model each of these three aspects in our  
110 ResearchAgent. Specifically, in order to imbibe re-  
111 lated work, the system begins with a core scientific  
112 paper and then explores a range of related papers  
113 through references and citation relationships. Fur-  
114 ther, to develop an encyclopedic view of related  
115 concepts, we build and then augment ResearchA-  
116 gent with an entity-centric knowledge store derived  
117 from co-occurrences of key concepts in the sci-  
118 entific literature. This repository is aimed at cap-

119 turing novel underlying relationships within and  
120 across domains, thereby increasing the chances of  
121 a cross-pollination of ideas. Finally, in order to sim-  
122 ulate robust feedback mechanisms, we instantiate  
123 a number of LLM-powered ReviewingAgents that  
124 help the ResearchAgent to iterate on research idea  
125 generation with constructive critiques. Crucially,  
126 these ReviewingAgents are prompted with evalua-  
127 tion criteria that are induced from real researchers’  
128 judgements, thus aligning them with actual scien-  
129 tific preferential standards. An illustration of our  
130 overarching system is provided in Figure 1.

131 We validate the effectiveness of ResearchAgent  
132 for research idea generation based on scientific liter-  
133 ature across multiple disciplines. Then, on a battery  
134 of tests conducted with both human- and model-  
135 based evaluations, we demonstrate that ResearchA-  
136 gent outperforms strong LLM-powered baselines  
137 by large margins, generating more clear, relevant,  
138 and significant ideas that are especially novel. Fur-  
139 thermore, analyses show the efficacy of our com-  
140 prehensive approach to modeling ResearchAgent:  
141 the entity-centric knowledge store and the itera-  
142 tive idea refinement steps help the system generate  
143 meaningfully better ideas compared with an instan-  
144 tiation that is purely based on prior related work.

145 These findings indicate the huge promise of AI-  
146 mediated research assistants, and our initial novel  
147 foray into scientific idea generation has important  
148 implications for future work that seeks to explore  
149 and improve upon the work we have proposed here.  
150 These include better support and operationalization  
151 to experimentally validate scientific ideas, and the  
152 design and evaluation of the utility of these systems  
153 to end users, applications, and industries.

## 154 2 Related Work

155 **Large Language Models** Large Language Mod-  
156 els (LLMs) have shown impressive performances  
157 across a wide range of tasks (OpenAI, 2023; Anil  
158 et al., 2023), including ones in advanced scientific  
159 fields such as mathematics, physics, medicine, and  
160 computer science (Romera-Paredes et al., 2023;  
161 Bran et al., 2023; Huang et al., 2023). A recent  
162 study on GPT-4 shows that it is capable of under-  
163 standing DNA sequences, designing biomolecules,  
164 predicting the behavior of molecular systems, and  
165 solving Partial Differential Equation (PDE) prob-  
166 lems (AI4Science and Quantum, 2023). However,  
167 LLMs have mainly been used for accelerating the  
168 experimental validation of already identified re-  
169 search ideas, 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, all these approaches perform idea generation in a localized manner and are designed to identify potential relationships between two variables or to generate textual descriptions about them, which may be sub-optimal to capture the complexity and multifaceted nature of real-world problems. Meanwhile, we target a significantly more challenging open-ended scenario, aiming to generate research ideas that involve not only generating novel research hypotheses, but also outlines of methods, and experimental designs for these hypotheses.

**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 – while similar – has the entirely different objective of narrowing the focus of LLMs to entities already present in an LLM’s context.

**Iterative Refinements with LLMs** Similar to humans, LLMs do not always generate optimal out-

puts on their first attempt. Drawing inspiration from humans who can iteratively refine their thoughts based on critiques from themselves and their peers, many recent studies (including some hypothesis generation work) have investigated the potential of LLMs to correct and refine their outputs, demonstrating that 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).

### 3 Method

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

#### 3.1 LLM-Powered Research Idea Generation

We begin by formally introducing the new problem of research idea generation, followed by an explanation of how LLMs are utilized 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. While we acknowledge that the real process by which humans conduct research is varied and complex to an extent well beyond the scope of this scientific study, we attempt to model a simulacrum in three systematic steps that would likely be maximally beneficial to a researcher seeking assistance from an AI system. These are namely, identifying novel research ideas, proposing methods to validate these ideas, and designing experiments to measure the success of these methods in relation to the ideas.

To accomplish the aforementioned steps, we utilize the existing literature (e.g., academic publications) as a primary source, which provides insights about existing knowledge along with gaps and unanswered questions<sup>1</sup>. 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 opera-

<sup>1</sup>We focus on the existing literature-based idea generation by following the paradigm that a *new idea* is more often than not just a new combination of old elements (Young, 2003).

267 tionalize  $f$  with LLMs, leveraging their capability  
268 to understand and generate academic text.

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

### 280 3.2 Knowledge-Augmented LLMs 281 for Research Idea Generation

282 We now turn to our primary focus of automati-  
283 cally generating research ideas with LLMs. Re-  
284 call that we aim to produce a complete idea con-  
285 sisting of the problem, method, and experiment  
286 design ( $o = [p, m, d]$ ), while using the exist-  
287 ing literature  $\mathcal{L}$  as a primary source of informa-  
288 tion. We operationalize this with LLMs by instan-  
289 tiating the aforementioned research idea genera-  
290 tion function  $f$  with LLM coupled with the task-  
291 specific template. Formally,  $p = \text{LLM}(\mathcal{T}_p(\mathcal{L}))$   
292 indicates the problem identification step, followed  
293 by  $m = \text{LLM}(\mathcal{T}_m(p, \mathcal{L}))$  for method development  
294 and  $d = \text{LLM}(\mathcal{T}_e(p, m, \mathcal{L}))$  for experiment design,  
295 which constitutes the full idea:  $o = [p, m, d]$ .

296 Following this general formulation, the impor-  
297 tant question to answer is how the body of scientific  
298 literature is leveraged for actually generating re-  
299 search ideas with LLMs. Here, we outline three key  
300 desiderata that contribute to the success of human  
301 researchers ideating novel research ideas: a broad  
302 and deep understanding of related work, an ency-  
303 clopedic perspective on the interconnectedness of  
304 concepts within and across scientific domains, and  
305 a community of peers who help iteratively improve  
306 ideas through constructive critiques. We describe  
307 our operationalization of these three desiderata us-  
308 ing the prior literature and LLMs in what follows.

309 **Citation Graph based Literature Survey** Due  
310 to the constraints on their input lengths and their  
311 reasoning abilities, particularly over very long con-  
312 texts (Liu et al., 2023), it is not possible to incorpo-  
313 rate all the existing publications from the literature  
314  $\mathcal{L}$  into the LLM input. Instead, we need to find a  
315 meaningful subset relevant to the problem at hand.  
316 To achieve this, we mirror the process followed by

317 human researchers, who expand their knowledge of  
318 a paper by perusing other papers that either cite or  
319 are cited by it. Concretely, for the LLM, we initiate  
320 its literature review process by providing a core  
321 paper  $l_0$  from  $\mathcal{L}$  and then selectively incorporat-  
322 ing subsequent papers  $\{l_1, \dots, l_n\}$  that are directly  
323 connected based on a citation graph. This proce-  
324 dure makes the LLM input for idea generation more  
325 manageable and coherent. In addition, we oper-  
326 ationalize the selection process of the core paper  
327 and its relevant citations with two design choices:  
328 1) the core paper is selected based on its citation  
329 count (e.g., exceeding 100 over 3 months) typi-  
330 cally indicating high impact; 2) its relevant papers  
331 (which may be potentially numerous) are further  
332 narrow-downed based on their similarities of ab-  
333 stracts with the core paper, ensuring a more focused  
334 and relevant set of related work.

335 **Entity-Centric Knowledge Augmentation** In  
336 order to model an encyclopedic view of inter-  
337 connected concepts, we must effectively design a  
338 framework to extract, store and effectively leverage  
339 the vast amount of knowledge in scientific litera-  
340 ture  $\mathcal{L}$ . In this work, we view entities as the atomic  
341 units of knowledge, which allows for ease of repre-  
342 sentation and accumulation over papers in a unified  
343 manner across different disciplines. For example,  
344 we can easily extract the term “database” whenever  
345 it appears in any paper, using existing off-the-shelf  
346 entity linking methods and then aggregate their  
347 linked occurrences into a knowledge store. Then, if  
348 the term “database” is prevalent within the realm of  
349 medical science but less so in hematology (which  
350 is a subdomain of medical science), the constructed  
351 knowledge store can capture the affinity between  
352 those two domains based on overlapping entities.  
353 This representational paradigm can then be used  
354 to suggest the term “database” when formulating  
355 the ideas about hematology. In other words, this  
356 approach enables providing novel and interdis-  
357 ciplinary insights by leveraging the interconnected-  
358 ness of entities across various fields.

359 Formally, we design the knowledge store as a  
360 two-dimensional matrix  $\mathcal{K} \in \mathcal{R}^{m \times m}$  where  $m$  is  
361 the total number of unique entities identified and  
362  $\mathcal{K}$  is implemented in a sparse format. This knowl-  
363 edge store is constructed by extracting entities over  
364 all the available scientific articles in literature  $\mathcal{L}^2$ ,  
365 which not only counts the co-occurrences between

<sup>2</sup>As extracting entities on all articles is computationally infeasible, we target papers appearing after May 01, 2023.

entity pairs within individual papers but also quantifies the count for each entity. Our approach is versatile, thus, we can use any entity linker (Wu et al., 2020). Also, despite the lack of entity linkers customized for the scientific domain, the off-the-shelf system proved capable of extracting key scientific entities, as shown in Table 15. Specifically, this linker tags and canonicalizes entities in a 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  $l$ <sup>3</sup>. 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  $e \in \mathcal{E}$ .

Given this knowledge store  $\mathcal{K}$ , our next goal is to enhance the previous vanilla research idea generation process implemented 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, 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 matrix  $\mathcal{K}$ , suitably normalized. We note that the formulation in Equation 2 is only one instance of operationalizing retrieval; this could be replaced with other retrieval strategies – for example, the embedding-based retrieval (discussions and results are provided in Appendix B.2). Hereafter, the instantiation of research proposal generation augmented with relevant entity-centric knowledge is formalized as follows:  $\mathbf{o} =$

$\text{LLM}(\mathcal{T}(\{l_0, \dots, l_n\}, \text{Ret}(\{l_0, \dots, l_n\}; \mathcal{K})))$ <sup>4</sup>. We call this knowledge-augmented LLM-powered idea generation approach `ResearchAgent`, and provide the templates to instantiate it in Tables 5, 6, and 7.

**Iterative Research Idea Refinements** Finally, in order to model a community of peers for idea improvement, we propose a set of LLM-powered reviewing agents (called `ReviewingAgents`). These agents provide the `ResearchAgent` with reviews and feedback according to specific criteria in order to help it iteratively improve idea generation.

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 8, 9, and 10). Then, with `ReviewingAgents`, each of the generated research ideas (problem, method, and experiment design) is separately evaluated according to its own specific five criteria<sup>5</sup>, which are provided in labels of Figure 2 and detailed in Table 11. Based on the reviews and feedback from `ReviewingAgents`, the `ResearchAgent` iteratively updates and refines its generation of 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 complex research ideas may not be aligned with the judgments of real human researchers. 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 made 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 research ideas and their associated scores (on a 5-point Likert scale annotated by human researchers having at least 3 papers) on every evaluation criterion. Then, we prompt the LLM with these human-annotated pairs and ask it to induce detailed descriptions for evaluation criteria (Lin et al., 2024) (See Tables 12, 13, and 14) that reflect the human preferences, which are then used as evaluation criteria by the `ReviewingAgents` in the evaluation prompt template  $\mathcal{T}$ .

<sup>4</sup>There may be additional knowledge sources (beyond the existing literature and entities) for research idea generation, and we leave exploring them as future work.

<sup>5</sup>We select the top five criteria which we consider as the most important, and leave exploring others as future work.

<sup>3</sup>Due to the extensive length of scientific publications, the target of entity extraction is restricted to titles and abstracts.

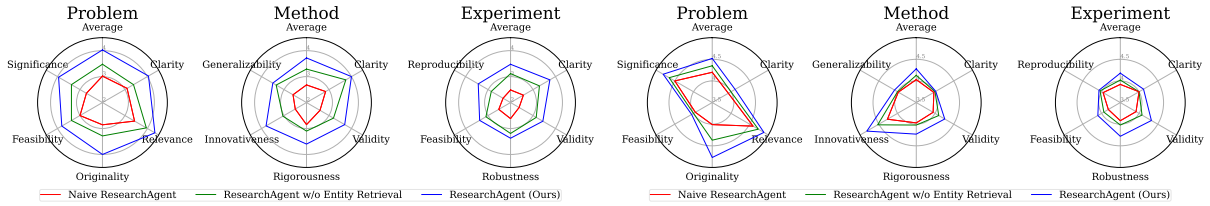


Figure 2: Main results on our research idea generation task with human- (left) and model-based (right) evaluations, where we report the score of each idea (problem, method, or experiment design) based on its own five criteria and their average score.

## 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 API<sup>6</sup>. From this, we select papers appearing after May 01, 2023, because LLMs that we use in our experiments are trained on data from the open web available before this point. This follows the procedure of existing literature-based hypothesis generation work (Qi et al., 2023). Then, we select high-impact papers (that have more than 20 citations) as core papers, mirroring human researchers’ tendency to leverage influential work, to ensure the high quality of the generated ideas. The resulting data is still very large; thus, we further randomly sample a subset of 300 papers as core papers to obtain a reasonably sized benchmark dataset. 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.

### 4.2 Baselines and Our Model

In this work, we target the novel task of research idea generation, for which there are no existing baselines that would serve as 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 Setup

Given our formulation of idea generation (Sec 3.1), there are no ground-truth answers to measure the

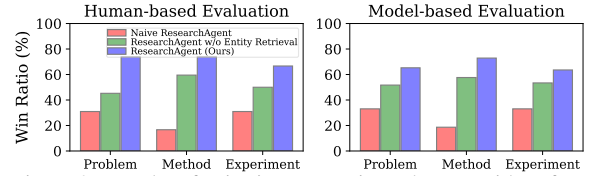


Figure 3: Results of pairwise comparisons between ideas from two of any different approaches, where we report the win ratio.

quality of the generated ideas. Meanwhile, exhaustively listing pairs of core papers and reference research ideas is suboptimal, 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. Thus, 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 and see Table 11 for their detailed descriptions). We ask the LLM-based 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. We provide the 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. As the generated ideas are knowledge-intensive, we carefully select annotators who are well-versed in the field and provide them with ideas that are highly relevant to their field of expertise. Specifically, we choose ten expert researchers who have authored at least three papers and ask them to judge only the ideas that are generated based on their own papers.

### 4.4 Implementation Details

We mainly use the GPT-4 (OpenAI, 2023) release from Nov 06 as the basis for all models, which is,

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

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

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

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), with papers from May 01, 2023, to Dec 31, 2023 (available from Semantic Scholar API) along with their references, which number 50,091 in total. We provide detailed 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 all the problems, methods, and experiment designs generated (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.

**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

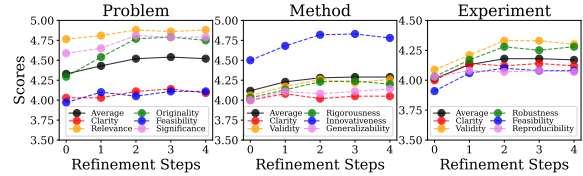


Figure 4: Results with varying the number of refinement steps.

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

Methods	Problem	Method	Experiment
ResearchAgent	<b>4.52</b>	<b>4.28</b>	<b>4.18</b>
- 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

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 full ResearchAgent is augmented with two different knowledge sources, namely relevant references and entities. To see their individual contribution, we perform an ablation study by either excluding one of the knowledge sources or replacing it with random elements. As shown in Table 2, each knowledge source appears to contribute to performance improvement, and the relevant references are especially helpful. We also note that providing random elements is more helpful than providing no elements at all; we hypothesize that this may be due to the LLM’s capability to filter out noise while still gaining incidental value from random inputs.

**Analysis on Human Alignment for Evaluation**

Recall that to align judgments from model-based

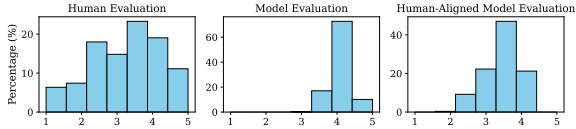


Figure 5: Distributions of model-based evaluation results with and without the human-induced score criteria alignment (middle and right), as well as human evaluation results (left).

616 evaluations with actual human preferences, we gener-  
 617 erated the evaluation criteria based on human eval-  
 618 uation results and used them as the criteria for  
 619 model-based evaluations. Figure 5 demonstrates  
 620 the efficacy of this strategy, presenting the score  
 621 distribution of human evaluation compared with the  
 622 distributions of model-based evaluations with and  
 623 without human alignment. We find that the score  
 624 distribution of model-based evaluations without  
 625 alignment is skewed and different from the score  
 626 distribution of human judgments. Meanwhile, after  
 627 aligning the model-based evaluations with human-  
 628 induced score criteria, the calibrated distribution  
 629 more closely resembles the distribution of humans.

630 **Correlation on Citation Counts** We further in-  
 631 vestigate whether a high-impact paper (when used  
 632 as a core paper) leads to high-quality research ideas.  
 633 To measure this, we bucketize all papers into three  
 634 groups by the number of their citations (using it as  
 635 a proxy for impact), and visualize the average score  
 636 of each bucket (with model-based evaluations) in  
 637 Figure 6. We observe that the research ideas gener-  
 638 ated from high-impact papers are generally of high  
 639 quality. Additionally, based on the paper distribu-  
 640 tion (See Figure 7) and for the ease of manual qual-  
 641 ity check, evaluation criteria for model-based eval-  
 642 uations are induced mainly with computer science  
 643 papers. To see whether those criteria are applica-  
 644 ble to diverse fields, we also compare a correlation  
 645 between scores of computer science papers and all  
 646 papers in Figure 6. From this, we observe that the  
 647 scores increase when the citation increases for both  
 648 domains, which may support the generalizability  
 649 of human-preference-induced evaluation criteria.

650 **Analysis using Different LLMs** To see how the  
 651 performance of ResearchAgent changes if an LLM  
 652 other than the GPT-4 is used, we conduct an auxil-  
 653 iary analysis instantiating the ResearchAgent with  
 654 different LLMs, such as Llama3, Mixtral, Qwen1.5,  
 655 and GPT-3.5 (Bai et al., 2023; Jiang et al., 2024),  
 656 and present the model-based evaluation results in  
 657 Table 3. From this, we find that the performance  
 658 with less capable models (other than GPT-4) drops  
 659 significantly, justifying our choice to not consider

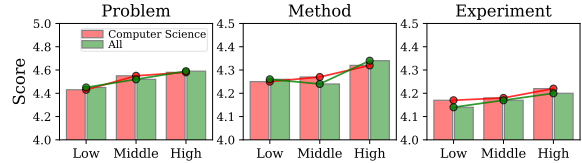


Figure 6: Results with bucketing papers based on citations.

Table 3: Results with different, open and proprietary LLMs.

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
Llama3 (8B)	Naive ResearchAgent	3.76	3.69	3.54
	ResearchAgent (Ours)	4.18	4.03	3.95
Mixtral (8x7B)	Naive ResearchAgent	3.31	3.27	3.20
	ResearchAgent (Ours)	3.28	3.35	3.31
Qwen1.5 (32B)	Naive ResearchAgent	3.64	3.74	3.66
	ResearchAgent (Ours)	4.02	3.97	3.94

660 weaker LLMs than GPT-4. Additionally, the perfor-  
 661 mance differences between the Naive ResearchA-  
 662 gent without knowledge augmentation and the full  
 663 ResearchAgent become marginal, for Mixtral and  
 664 GPT-3.5, which indicates that they might simply  
 665 not be capable of capturing complex concepts and  
 666 their relationships across different scientific papers.  
 667 This is unsurprising if taken in the context of the  
 668 emergent abilities of LLMs for complex reasoning  
 669 (but not in smaller LMs) (Wei et al., 2022).

## 6 Conclusion

670 In this work, we presented ResearchAgent – a sys-  
 671 tem that aims to assist researchers in their workflow  
 672 by automatically generating research ideas, which  
 673 consists of novel problem identification, method de-  
 674 velopment, and experiment design. Drawing inspi-  
 675 ration from the human process of research ideation,  
 676 we developed an approach that simultaneously con-  
 677 ducts a broad and deep review of relevant literature,  
 678 leverages encyclopedic knowledge through inter-  
 679 connected concepts across domains to help cross-  
 680 pollination of ideas, and leverages a community of  
 681 reviewing agents to provide constructive critiques  
 682 for iteratively refining the research ideas. Through  
 683 human and model-based evaluations, we showed  
 684 that ResearchAgent generates ideas that are more  
 685 creative, valid, and clear than ones from baselines.  
 686 While we envision ResearchAgent as a collabora-  
 687 tive partner for scientists, this initial foray has  
 688 only demonstrated early signs of the promise of  
 689 AI-mediated research assistants. There are multi-  
 690 ple important avenues of future research to pursue,  
 691 including improving and building upon ResearchA-  
 692 gent, operationalizing experimental validation of its  
 693 research hypotheses, and measuring the real-world  
 694 value it brings to researchers and their productivity.  
 695



## 696 Limitations

697 ResearchAgent has some clear limitations that we  
698 hope to address in future work. First, recall that we  
699 built the entity-centric knowledge store to propose  
700 beneficial entities during idea generation; how-  
701 ever this store is constructed by extracting entities  
702 from the titles and abstracts of a limited number of  
703 publications (due to the costs of processing them)  
704 thereby precluding a large number of other entities  
705 and their interconnectedness. In addition, the num-  
706 ber of entities that we obtain from the BLINK entity  
707 linker (Wu et al., 2020) amounts to 3 per paper on  
708 average, indicating limited coverage (it is an open-  
709 domain linker after all). We argue that to build  
710 a more comprehensive entity-centric knowledge  
711 store, future work will not only need to extend the  
712 content (including the main texts of publications)  
713 and the volume of papers for entity extraction, but  
714 also improve the capability of the entity linker itself  
715 to more accurately extract scientific terms within  
716 the literature. Lastly, since our ResearchAgent is  
717 powered by LLMs, similar to any other approaches  
718 based on LLMs, it may hallucinate the generated  
719 research ideas. While our proposed ResearchAgent  
720 can partially mitigate this problem by augmenting  
721 LLMs with additional elements, such as references  
722 to the target paper and greater entity-centric knowl-  
723 edge, which help ground the generation process in  
724 more accurate and relevant information, validating  
725 these generated research ideas is essential to truly  
726 accelerate scientific research.

## 727 Ethics Statement

728 We are aware that the ResearchAgent may have the  
729 potential to be misused for nefarious purposes, such  
730 as generating research ideas about new explosives,  
731 malicious software, and invasive surveillance tools.  
732 Notably, this vulnerability is not unique to our ap-  
733 proach but a common challenge faced by existing  
734 LLMs that possess significant creative and reason-  
735 ing capabilities, occasionally generating content  
736 that may be deemed undesirable. Consequently, it  
737 underscores the necessity to enhance the robustness  
738 and safety of LLMs more broadly.

## 739 References

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

Rohan Anil, Sebastian Borgeaud, Yonghui Wu, Jean-  
Baptiste Alayrac, Jiahui Yu, Radu Soricut, Johan  
Schalkwyk, Andrew M. Dai, Anja Hauth, Katie Mil-  
lican, David Silver, Slav Petrov, Melvin Johnson,  
Ioannis Antonoglou, Julian Schrittwieser, Amelia  
Glaese, Jilin Chen, Emily Pitler, Timothy P. Lilli-  
crap, Angeliki Lazaridou, Orhan Firat, James Molloy,  
Michael Isard, Paul Ronald Barham, Tom Hennig-  
an, Benjamin Lee, Fabio Viola, Malcolm Reynolds,  
Yuanzhong Xu, Ryan Doherty, Eli Collins, Clemens  
Meyer, Eliza Rutherford, Erica Moreira, Kareem  
Ayoub, Megha Goel, George Tucker, Enrique Pi-  
queras, Maxim Krikun, Iain Barr, Nikolay Savinov,  
Ivo Danihelka, Becca Roelofs, Anaïs White, Anders  
Andreassen, Tamara von Glehn, Lakshman Yagati,  
Mehran Kazemi, Lucas Gonzalez, Misha Khalman,  
Jakub Sygnowski, and et al. 2023. [Gemini: A family  
of highly capable multimodal models](#). *arXiv preprint  
arXiv:2312.11805*. 744-762

Jinheon Baek, Alham Fikri Aji, and Amir Saffari. 2023. [Knowledge-augmented language model prompting for zero-shot knowledge graph question answering](#). In *Proceedings of the 1st Workshop on Natural Language Reasoning and Structured Explanations (NLRSE)*, pages 78–106, Toronto, Canada. Association for Computational Linguistics. 763-769

Jinheon Baek, Nirupama Chandrasekaran, Silviu Cucerzan, Allen Herring, and Sujay Kumar Jauhar. 2024. [Knowledge-augmented large language models for personalized contextual query suggestion](#). *WWW*. 770-773

Jinze Bai, Shuai Bai, Yunfei Chu, Zeyu Cui, Kai Dang, Xiaodong Deng, Yang Fan, Wenheng Ge, Yu Han, Fei Huang, Binyuan Hui, Luo Ji, Mei Li, Junyang Lin, Runji Lin, Dayiheng Liu, Gao Liu, Chengqiang Lu, K. Lu, Jianxin Ma, Rui Men, Xingzhang Ren, Xuancheng Ren, Chuanqi Tan, Sinan Tan, Jianhong Tu, Peng Wang, Shijie Wang, Wei Wang, Shengguang Wu, Benfeng Xu, Jin Xu, An Yang, Hao Yang, Jian Yang, Jian Yang, Shusheng Yang, Yang Yao, Bowen Yu, Yu Bowen, Hongyi Yuan, Zheng Yuan, Jianwei Zhang, Xing Zhang, Yichang Zhang, Zhenru Zhang, Chang Zhou, Jingren Zhou, Xiaohuan Zhou, and Tianhang Zhu. 2023. [Qwen technical report](#). *arXiv preprint arXiv:2309.16609*. 774-787

Andrés M Bran, Sam Cox, Oliver Schilter, Carlo Baldasari, Andrew D. White, and Philippe Schwaller. 2023. [Chemcrow: Augmenting large-language models with chemistry tools](#). 788-791

Jacob Cohen. 1960. [A coefficient of agreement for nominal scales](#). *Educational and Psychological Measurement*, 20:37 – 46. 792-794

Michael Fire and Carlos Guestrin. 2019. [Over-optimization of academic publishing metrics: Observing goodhart’s law in action](#). *GigaScience*, 8. 795-797

Jinlan Fu, See-Kiong Ng, Zhengbao Jiang, and Pengfei Liu. 2023. [Gptscore: Evaluate as you desire](#). *arXiv preprint arXiv:2302.04166*. 798-800

801	Deep Ganguli, Amanda Askell, Nicholas Schiefer,	Aman Madaan, Niket Tandon, Prakhar Gupta, Skyler	858
802	Thomas I. Liao, Kamile Lukosiute, Anna Chen,	Hallinan, Luyu Gao, Sarah Wiegrefe, Uri Alon,	859
803	Anna Goldie, Azalia Mirhoseini, Catherine Olsson,	Nouha Dziri, Shrimai Prabhumoye, Yiming Yang,	860
804	Danny Hernandez, Dawn Drain, Dustin Li, Eli Tran-	Shashank Gupta, Bodhisattwa Prasad Majumder,	861
805	Johnson, Ethan Perez, Jackson Kernion, Jamie Kerr,	Katherine Hermann, Sean Welleck, Amir Yazdan-	862
806	Jared Mueller, Joshua Landau, Kamal Ndousse, Ka-	bakhsh, and Peter Clark. 2023. <a href="#">Self-refine: Iterative refinement with self-feedback</a> . 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> .	863
807	rina Nguyen, Liane Lovitt, Michael Sellitto, Nelson		864
808	Elhage, Noemí Mercado, Nova DasSarma, Oliver		865
809	Rausch, Robert Lasenby, Robin Larson, Sam Ringer,		866
810	Sandipan Kundu, Saurav Kadavath, Scott Johnston,		867
811	Shauna Kravec, Sheer El Showk, Tamera Lanham,		868
812	Timothy Telleen-Lawton, Tom Henighan, Tristan		
813	Hume, Yuntao Bai, Zac Hatfield-Dodds, Ben Mann,	R.K. Nadkarni, David Wadden, Iz Beltagy, Noah A.	869
814	Dario Amodei, Nicholas Joseph, Sam McCandlish,	Smith, Hannaneh Hajishirzi, and Tom Hope. 2021.	870
815	Tom Brown, Christopher Olah, Jack Clark, Samuel R.	<a href="#">Scientific language models for biomedical knowl-</a>	871
816	Bowman, and Jared Kaplan. 2023. <a href="#">The capacity</a>	<a href="#">edge base completion: An empirical study</a> . <i>ArXiv</i> ,	872
817	<a href="#">for moral self-correction in large language models</a> .	abs/2106.09700.	873
818	<a href="#">arXiv preprint arXiv:2302.07459</a> .		
819	Sam Henry and Bridget T. McInnes. 2017. <a href="#">Literature</a>	OpenAI. 2023. <a href="#">GPT-4 technical report</a> . <i>arXiv preprint</i>	874
820	<a href="#">based discovery: Models, methods, and trends</a> . <i>Journal</i>	<i>arXiv:2303.08774</i> .	875
821	<a href="#">of biomedical informatics</a> , 74:20–32.		
822	Tom Hope, Doug Downey, Daniel S. Weld, Oren Et-	W. Pirie. 2006. <i>Spearman Rank Correlation Coefficient</i> ,	876
823	zioni, and Eric Horvitz. 2023. <a href="#">A computational</a>	volume 8.	877
824	<a href="#">inflection for scientific discovery</a> . <i>Commun. ACM</i> ,		
825	66(8):62–73.	Biqing Qi, Kaiyan Zhang, Haoxiang Li, Kai Tian, Si-	878
826	Qian Huang, Jian Vora, Percy Liang, and Jure Leskovec.	hang Zeng, Zhang-Ren Chen, and Bowen Zhou. 2023.	879
827	2023. <a href="#">Benchmarking large language models as AI</a>	<a href="#">Large language models are zero shot hypothesis pro-</a>	880
828	<a href="#">research agents</a> . <i>arXiv preprint arXiv:2310.03302</i> .	<a href="#">posers</a> . <i>arXiv preprint arXiv:2311.05965</i> .	881
829	Albert Q. Jiang, Alexandre Sablayrolles, Antoine	Ori Ram, Yoav Levine, Itay Dalmedigos, Dor Muhlgay,	882
830	Roux, Arthur Mensch, Blanche Savary, Chris Bam-	Amnon Shashua, Kevin Leyton-Brown, and Yoav	883
831	ford, Devendra Singh Chaplot, Diego de Las Casas,	Shoham. 2023. <a href="#">In-context retrieval-augmented lan-</a>	884
832	Emma Bou Hanna, Florian Bressand, Gianna	<a href="#">guage models</a> . <i>Transactions of the Association for</i>	885
833	Lengyel, Guillaume Bour, Guillaume Lample,	<i>Computational Linguistics</i> , 11:1316–1331.	886
834	L’elio Renard Lavaud, Lucile Saulnier, Marie-	Bernardino Romera-Paredes, Mohammadamin	887
835	Anne Lachaux, Pierre Stock, Sandeep Subramanian,	Barekatain, Alexander Novikov, Matej Balog,	888
836	Sophia Yang, Szymon Antoniak, Teven Le Scao,	M Pawan Kumar, Emilien Dupont, Francisco J. R.	889
837	Théophile Gervet, Thibaut Lavril, Thomas Wang,	Ruiz, Jordan S. Ellenberg, Pengming Wang, Omar	890
838	Timothée Lacroix, and William El Sayed. 2024. <a href="#">Mix-</a>	Fawzi, Pushmeet Kohli, Alhussein Fawzi, Josh	891
839	<a href="#">tral of experts</a> . <i>arXiv preprint arXiv: 2401.04088</i> .	Grochow, Andrea Lodi, Jean-Baptiste Mouret, Talia	892
840	Angeliki Lazaridou, Elena Gribovskaya, Wojciech	Ringer, and Tao Yu. 2023. <a href="#">Mathematical discoveries</a>	893
841	Stokowiec, and Nikolai Grigorev. 2022. <a href="#">Internet-</a>	<a href="#">from program search with large language models</a> .	894
842	<a href="#">augmented language models through few-shot</a>	<i>Nature</i> , 625:468 – 475.	895
843	<a href="#">prompting for open-domain question answering</a> .	Weijia Shi, Sewon Min, Michihiro Yasunaga, Min-	896
844	<i>arXiv preprint arXiv:2203.05115</i> .	joon Seo, Rich James, Mike Lewis, Luke Zettle-	897
845	Ying-Chun Lin, Jennifer Neville, Jack W Stokes,	moyer, and Wen tau Yih. 2023. <a href="#">Replug: Retrieval-</a>	898
846	Longqi Yang, Tara Safavi, Mengting Wan, Scott	<a href="#">augmented black-box language models</a> . <i>arXiv</i>	899
847	Counts, Siddharth Suri, Reid Andersen, Xiaofeng	<i>preprint arXiv:2301.12652</i> .	900
848	Xu, Deepak Gupta, Sujay Kumar Jauhar, Xia Song,	Kumar Shridhar, Koustuv Sinha, Andrew Cohen, Tianlu	901
849	Georg Buscher, Saurabh Tiwary, Brent Hecht, and	Wang, Ping Yu, Ram Pasunuru, Mrinmaya Sachan,	902
850	Jaime Teevan. 2024. <a href="#">Interpretable user satisfaction</a>	Jason Weston, and Asli Celikyilmaz. 2023. <a href="#">The ART</a>	903
851	<a href="#">estimation for conversational systems with large lan-</a>	<a href="#">of LLM refinement: Ask, refine, and trust</a> . <i>arXiv</i>	904
852	<a href="#">guage models</a> . <i>arXiv preprint arXiv:2403.12388</i> .	<i>preprint arXiv:2311.07961</i> .	905
853	Nelson F. Liu, Kevin Lin, John Hewitt, Ashwin Paran-	Don R. Swanson. 1986. <a href="#">Undiscovered public knowl-</a>	906
854	jape, Michele Bevilacqua, Fabio Petroni, and Percy	<a href="#">edge</a> . <i>The Library Quarterly</i> , 56:103–118.	907
855	Liang. 2023. <a href="#">Lost in the middle: How language mod-</a>	Justin Sybrandt, Ilya Tyagin, M. Shtutman, and Ilya	908
856	<a href="#">els use long contexts</a> . <i>Transactions of the Association</i>	Safro. 2020. <a href="#">Agatha: Automatic graph mining and</a>	909
857	<a href="#">for Computational Linguistics</a> , 12:157–173.	<a href="#">transformer based hypothesis generation approach</a> .	910
		<i>Proceedings of the 29th ACM International Confer-</i>	911
		<i>ence on Information &amp; Knowledge Management</i> .	912

913	Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton-Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony Hartshorn, Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel Kloumann, Artem Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov, Pushkar Mishra, Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan Saladi, Alan Schelten, Ruan Silva, Eric Michael Smith, Ranjan Subramanian, Xiaoqing Ellen Tan, Binh Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zheng Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurélien Rodriguez, Robert Stojnic, Sergey Edunov, and Thomas Scialom. 2023. <a href="#">Llama 2: Open foundation and fine-tuned chat models</a> . <i>arXiv preprint arXiv:2307.09288</i> .	
936	Vahe Tshitoyan, John Dagdelen, Leigh Weston, Alex Dunn, Ziqin Rong, Olga Vitalievna Kononova, Kristin A. Persson, Gerbrand Ceder, and Anubhav Jain. 2019. <a href="#">Unsupervised word embeddings capture latent knowledge from materials science literature</a> . <i>Nature</i> , 571:95 – 98.	
942	Jessica Vamathevan, Dominic Clark, Paul Czodrowski, Ian Dunham, Edgardo Ferran, George Lee, Bin Li, Anant Madabhushi, Parantu Shah, Michaela Spitzer, and Shanrong Zhao. 2019. <a href="#">Applications of machine learning in drug discovery and development</a> . <i>Nature reviews. Drug discovery</i> , 18(6):463—477.	
948	Hanchen Wang, Tianfan Fu, Yuanqi Du, Wenhao Gao, Kexin Huang, Ziming Liu, Payal Chandak, Shengchao Liu, Peter Van Katwyk, Andreea Deac, Anima Anandkumar, Karianne Bergen, Carla P. Gomes, Shirley Ho, Pushmeet Kohli, Joan Lasenby, Jure Leskovec, Tie-Yan Liu, Arjun Manrai, Debora S. Marks, Bharath Ramsundar, Le Song, Jimeng Sun, Jian Tang, Petar Velickovic, Max Welling, Linfeng Zhang, Connor W. Coley, Yoshua Bengio, and Marinka Zitnik. 2023a. <a href="#">Scientific discovery in the age of artificial intelligence</a> . <i>Nat.</i> , 620(7972):47–60.	
959	Qingyun Wang, Doug Downey, Heng Ji, and Tom Hope. 2023b. <a href="#">Learning to generate novel scientific directions with contextualized literature-based discovery</a> . <i>arXiv preprint arXiv:2305.14259</i> .	
963	Jason Wei, Yi Tay, Rishi Bommasani, Colin Raffel, Barret Zoph, Sebastian Borgeaud, Dani Yogatama, Maarten Bosma, Denny Zhou, Donald Metzler, Ed Huai hsin Chi, Tatsunori Hashimoto, Oriol Vinyals, Percy Liang, Jeff Dean, and William Fedus. 2022. <a href="#">Emergent abilities of large language models</a> . <i>arXiv preprint arXiv:2206.07682</i> .	
970	Sean Welleck, Ximing Lu, Peter West, Faeze Brahman, Tianxiao Shen, Daniel Khoshabi, and Yejin	
	Choi. 2023. <a href="#">Generating sequences by learning to self-correct</a> . In <i>The Eleventh International Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023</i> . OpenReview.net.	972 973 974 975
	Ledell Wu, Fabio Petroni, Martin Josifoski, Sebastian Riedel, and Luke Zettlemoyer. 2020. <a href="#">Scalable zero-shot entity linking with dense entity retrieval</a> . In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)</i> , pages 6397–6407, Online. Association for Computational Linguistics.	976 977 978 979 980 981 982
	Yike Wu, Nan Hu, Sheng Bi, Guilin Qi, J. Ren, Anhuan Xie, and Wei Song. 2023. <a href="#">Retrieve-rewrite-answer: A kg-to-text enhanced llms framework for knowledge graph question answering</a> . <i>arXiv preprint arXiv:2309.11206</i> .	983 984 985 986 987
	Zonglin Yang, Xinya Du, Junxian Li, Jie Zheng, Soujanya Poria, and Erik Cambria. 2023. <a href="#">Large language models for automated open-domain scientific hypotheses discovery</a> . <i>arXiv preprint arXiv:2309.02726</i> .	988 989 990 991 992
	J. Young. 2003. <i>A Technique for Producing Ideas</i> . McGraw Hill LLC.	993 994
	Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric P. Xing, Haoteng Zhang, Joseph Gonzalez, and Ion Stoica. 2023. <a href="#">Judging llm-as-a-judge with mt-bench and chatbot arena</a> . <i>arXiv preprint arXiv:2306.05685</i> .	995 996 997 998 999 1000

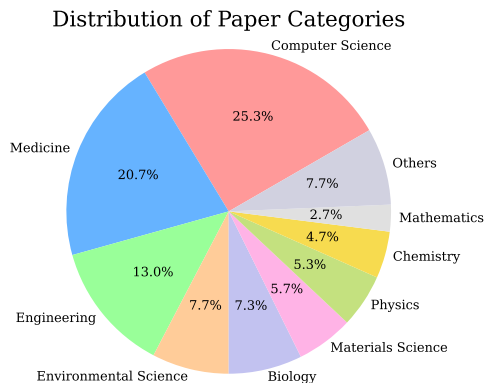


Figure 7: Visualization of the distribution of disciplines for all core papers, selected for research idea generation.

## 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 API<sup>7</sup>. 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 researchers from the United States and South Korea, majoring in computer science, medicine, and biology, each with a minimum of 3 published papers. For annotation, they were provided with a 6-page guideline document, which includes the task instruction and annotation examples. After reading this document, the annotators access the Label Studio platform, on which they first read the title and abstract of the target paper, and then review and evaluate the generated research ideas from different models. During the evaluation process, they are allowed to use any external tools, such as web searches. We note that they were compensated at a rate of \$22.20 per hour. Also, on average, for one hour, they evaluated 3 sets of research ideas (that are generated from their own papers), with each set comprising three sub-ideas (the 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

<sup>7</sup><https://www.semanticscholar.org/product/api>

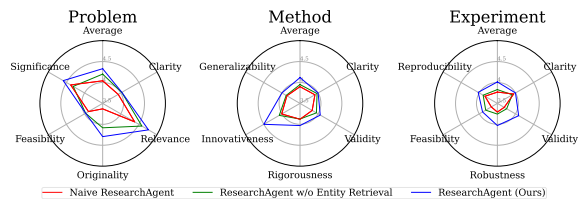


Figure 8: Results on our research idea generation task with model-based evaluation, where we exclude refinement steps.

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 5, Table 6, and Table 7, respectively.

### 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 8, Table 9, and Table 10, 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 12, Table 13, and Table 14, 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.

Table 4: Results with different entity retrieval strategies.

Methods	Problem	Method	Experiment
ResearchAgent			
- w/ Co-occurrence-based Retrieval	<b>4.52</b>	4.28	<b>4.18</b>
- w/ Embedding-based Retrieval	4.49	<b>4.34</b>	4.16
- w/o Entity Retrieval	4.35	4.13	4.02

## B.2 Analysis with Different Entity Retrieval

To see the effectiveness of different entity retrieval strategies, we perform additional experiments, replacing the co-occurrence-based entity retrieval in Equation 2 to the contextual embedding-based retrieval. Notably, this contextual embedding-based retrieval approach uses the entities that have the highest similarity to the entities appearing in the current literature (i.e., core paper and its references) used for idea generation, where the similarities are calculated based on embedding-level similarities between entities over the latent space represented by the entity linker (Wu et al., 2020). Therefore, unlike the previous co-occurrence-based entity retrieval that may retrieve entities that have opposite concepts to the main idea of the current core paper (since we often mention limitations of previous work along with the proposed ideas), this embedding-based approach may provide the ResearchAgent with mostly the entities having similar concepts to the core paper. Nevertheless, as shown in Table 4, the results with the strategy of entity co-occurrence-based retrieval are comparable to the results with the new embedding-based contextual retrieval. These results might confirm that there is not much difference in the quality of entity retrieval among those two strategies, i.e., most entities retrieved from the co-occurrence-based retrieval are contextually relevant for generating research ideas.

## B.3 Examples

We provide examples of generated research ideas (including problems, methods, and experiment designs) in Table 15.

Table 5: The prompt used in the full instantiation of ResearchAgent for problem identification.

Types	Texts
<b>System Message</b>	<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>
<b>User Message</b>	<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"> <li>- 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.</li> <li>- 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.</li> <li>- 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.</li> </ul> <p>Your approach should be systematic:</p> <ul style="list-style-type: none"> <li>- Start by thoroughly reading the title and abstract of the target paper to understand its core focus.</li> <li>- 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.</li> <li>- 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.</li> </ul> <p>I am going to provide the target paper, related papers, and entities, as follows:            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 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']}            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            Problem:            Rationale:</p>

Table 6: The prompt used in the full instantiation of ResearchAgent for method development.

Types	Texts
<b>System Message</b>	<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>
<b>User Message</b>	<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"> <li>- 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.</li> <li>- 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.</li> <li>- 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.</li> </ul> <p>Your approach should be systematic:</p> <ul style="list-style-type: none"> <li>- Start by thoroughly reading the research problem and its rationale, to understand your primary focus.</li> <li>- Next, proceed to review the titles and abstracts of existing studies, to gain a broader perspective and insights relevant to the primary research topic.</li> <li>- 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.</li> </ul> <p>I am going to provide the research problem, existing studies (target paper &amp; 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 7: 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"> <li>- The research problem has been formulated based on an in-depth review of existing studies and a potential exploration of relevant entities.</li> <li>- The scientific method has been proposed to tackle the research problem, which has been informed by insights gained from existing studies and relevant entities.</li> <li>- 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.</li> <li>- 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.</li> </ul> <p>Your approach should be systematic:</p> <ul style="list-style-type: none"> <li>- Start by thoroughly reading the research problem and its rationale followed by the proposed method and its rationale, to pinpoint your primary focus.</li> <li>- Next, proceed to review the titles and abstracts of existing studies, to gain a broader perspective and insights relevant to the primary research topic.</li> <li>- 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.</li> </ul>
User Message	<p>I am going to provide the research problem, scientific method, existing studies (target paper &amp; 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 8: The prompt used in the full instantiation of ReviewingAgent for problem validation.

Types	Texts
<b>System Message</b>	<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>
<b>User Message</b>	<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"> <li>- 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.</li> </ul> <p>The existing studies (target paper &amp; 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"> <li>- Start by thoroughly reading the research problem and its rationale, keeping in mind the context provided by the existing studies mentioned above.</li> <li>- Next, generate a review and feedback that should be constructive, helpful, and concise, focusing on the {metric} of the problem.</li> <li>- 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}</li> </ul> <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 9: The prompt used in the full instantiation of ReviewingAgent for method validation.

Types	Texts
<b>System Message</b>	<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>
<b>User Message</b>	<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"> <li>- 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.</li> <li>- 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.</li> </ul> <p>The research problem and existing studies (target paper &amp; 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"> <li>- 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.</li> <li>- Next, generate a review and feedback that should be constructive, helpful, and concise, focusing on the {metric} of the method.</li> <li>- 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}</li> </ul> <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 10: The prompt used in the full instantiation of ReviewingAgent for experiment design validation.

Types	Texts
<b>System Message</b>	<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>
<b>User Message</b>	<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"> <li>- The research problem has been formulated based on an in-depth review of existing studies and a potential exploration of relevant entities.</li> <li>- The scientific method has been proposed to tackle the research problem, which has been informed by insights gained from existing studies and relevant entities.</li> <li>- 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.</li> </ul> <p>The research problem, scientific method, and existing studies (target paper &amp; 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"> <li>- 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.</li> <li>- Next, generate a review and feedback that should be constructive, helpful, and concise, focusing on the {metric} of the experiment.</li> <li>- 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}</li> </ul> <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 11: The criteria used for evaluating research ideas: problems, methods, and experiment designs.

Types	Criteria	Texts
<b>Problem</b>	<b>Clarity</b>	It assesses whether the problem is defined in a clear, precise, and understandable manner.
	<b>Relevance</b>	It measures whether the problem is pertinent and applicable to the current field or context of study.
	<b>Originality</b>	It evaluates whether the problem presents a novel challenge or unique perspective that has not been extensively explored before.
	<b>Feasibility</b>	It examines whether the problem can realistically be investigated or solved with the available resources and within reasonable constraints.
	<b>Significance</b>	It assesses the importance and potential impact of solving the problem, including its contribution to the field or its broader implications.
<b>Method</b>	<b>Clarity</b>	It assesses whether the method is described in a clear, precise, and understandable manner that allows for replication and comprehension of the approach.
	<b>Validity</b>	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.
	<b>Rigorousness</b>	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.
	<b>Innovativeness</b>	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.
	<b>Generalizability</b>	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.
<b>Experiment</b>	<b>Clarity</b>	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.
	<b>Validity</b>	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.
	<b>Robustness</b>	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.
	<b>Feasibility</b>	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.
	<b>Reproducibility</b>	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 12: 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
<b>Problem</b>	<b>Clarity</b>	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.
<b>Relevance</b>	<b>Relevance</b>	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.
<b>Originality</b>	<b>Originality</b>	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.
<b>Feasibility</b>	<b>Feasibility</b>	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.
<b>Significance</b>	<b>Significance</b>	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 13: 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"> <li>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.</li> <li>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.</li> <li>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.</li> <li>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.</li> <li>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.</li> </ol>
		Validity
	Rigorousness	
		Innovativeness
	Generalizability	

Table 14: 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
	<b>Clarity</b>	<ol style="list-style-type: none"> <li>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.</li> <li>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.</li> <li>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.</li> <li>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.</li> <li>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.</li> </ol>
	<b>Validity</b>	<ol style="list-style-type: none"> <li>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.</li> <li>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.</li> <li>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.</li> <li>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.</li> <li>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.</li> </ol>
<b>Experiment</b>	<b>Robustness</b>	<ol style="list-style-type: none"> <li>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.</li> <li>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.</li> <li>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.</li> <li>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.</li> <li>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.</li> </ol>
	<b>Feasibility</b>	<ol style="list-style-type: none"> <li>1. The experiment design is fundamentally unfeasible, with insurmountable resource, time, or technological constraints that make implementation virtually impossible within the proposed framework.</li> <li>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.</li> <li>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.</li> <li>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.</li> <li>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.</li> </ol>
	<b>Reproducibility</b>	<ol style="list-style-type: none"> <li>1. The experiment design lacks critical details, making it virtually impossible for other researchers to replicate the study under the same conditions or methodologies.</li> <li>2. The experiment provides some essential information but omits significant details needed for replication, leading to considerable ambiguity in methodology or conditions.</li> <li>3. The experiment design includes sufficient details for replication, but lacks clarity or completeness in certain areas, posing challenges for seamless reproducibility.</li> <li>4. The experiment is well-documented with clear, detailed instructions and methodologies that allow for consistent replication, albeit with minor areas for improvement.</li> <li>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.</li> </ol>

Table 15: The examples of research idea generation results from the proposed full ResearchAgent.

Index	Types	Texts
1	Input	<p><b>Title:</b> Knowledge-Augmented Language Model Prompting for Zero-Shot Knowledge Graph Question Answering</p> <p><b>Abstract:</b> 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><b>Entities:</b> ['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	<p>Developing a Multimodal Knowledge-Aware Prompting System for Multilingual Zero-Shot Question Answering across Structured and Unstructured Data Sources</p> <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>
	Method	<ol style="list-style-type: none"> <li>1. Interconnected Multimodal Knowledge Orchestration (IMKO): <ul style="list-style-type: none"> <li>- 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.</li> <li>- 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.</li> </ul> </li> <li>2. Linguistically Inclusive Retrieval Engine (LIRE): <ul style="list-style-type: none"> <li>- LIRE extends EMKA with an emphasis on semantic understanding, using transformer-based models trained on diverse datasets, including idiomatic and cultural nuances across languages.</li> <li>- 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.</li> </ul> </li> <li>3. Prompt Learning and Optimization Nexus (PLON): <ul style="list-style-type: none"> <li>- Building on AMPL, PLON focuses on creating a library of optimized prompts categorized by linguistic features and data modalities, using Bayesian optimization algorithms.</li> <li>- 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.</li> </ul> </li> <li>4. Cross-Modal Integrative Learning System (C-MILS): <ul style="list-style-type: none"> <li>- 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.</li> <li>- 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.</li> </ul> </li> <li>5. User-Centric Adaptation and Privacy Framework (UCAPF): <ul style="list-style-type: none"> <li>- 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.</li> <li>- Aligns with GDPR and other privacy frameworks to fortify trust and address the ethical use of data in dynamic learning environments.</li> </ul> </li> <li>6. Global Evaluation and Reflective Testing System (GERTS): <ul style="list-style-type: none"> <li>- 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.</li> <li>- Offers a structured approach to cross-cultural evaluation, including the use of demographic and regional diversity in forming test cohorts.</li> </ul> </li> </ol>

*Continued on the next page*



Table 15 – Continued from the previous page

Index	Types	Texts
		<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"> <li>1. Detailed System Implementation Plan: <ul style="list-style-type: none"> <li>- Provide a publicly accessible project roadmap with specific milestones, resource allocation, and timelines.</li> </ul> </li> <li>2. Dataset Curation with Clear Guidelines: <ul style="list-style-type: none"> <li>- Publish precise annotation guidelines with strategies to prevent bias.</li> <li>- Document the dataset assembly process, including source selection and data processing procedures.</li> </ul> </li> <li>3. Transparent System Training: <ul style="list-style-type: none"> <li>- Offer a detailed training protocol with hyperparameters, optimization strategies, and Bayesian optimization processes used in PLON.</li> </ul> </li> <li>4. Structured Zero-Shot Evaluation: <ul style="list-style-type: none"> <li>- Outline evaluation metrics derived from GERTS with benchmark datasets to test zero-shot capabilities.</li> </ul> </li> </ol>
	<b>Experiment</b>	<ol style="list-style-type: none"> <li>5. Clearer Interdisciplinary Evaluation Protocol: <ul style="list-style-type: none"> <li>- Specify the composition of the evaluation committee, criteria for assessments, and methods for integrating the feedback.</li> </ul> </li> <li>6. Iterative Improvement with Validation Metrics: <ul style="list-style-type: none"> <li>- Describe statistical methods for reflective assessment and continuous improvement, aligned with multi-tiered GERTS methodology.</li> </ul> </li> <li>7. User-Centric Design and Privacy Compliance Evaluation: <ul style="list-style-type: none"> <li>- Structure user studies with targeted data points to assess usability and cultural adaptability.</li> <li>- Outline privacy compliance protocols to adhere to international standards.</li> </ul> </li> <li>8. Detailed Global Scalability Evaluation Method: <ul style="list-style-type: none"> <li>- Define evaluation metrics for scalability tests and describe diverse infrastructural setups.</li> </ul> </li> <li>9. Enhanced Reporting for Reproducibility: <ul style="list-style-type: none"> <li>- Commit to creating a comprehensive report with precise specifications, configurations, and instructions for replication purposes.</li> <li>- Utilize GitHub for version-controlled deposition of code and datasets, and arXiv for openly accessible experiment protocols and findings.</li> </ul> </li> </ol>
	<b>Input</b>	<p><b>Title:</b> Test-Time Self-Adaptive Small Language Models for Question Answering</p> <p><b>Abstract:</b> 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><b>Entities:</b> ['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>
	<b>Problem</b>	Developing a Scalable, Domain-Adaptive Test-Time Training Protocol for Low-Resource Language QA Using Small Language Models

Continued on the next page

**Table 15 – Continued from the previous page**

Index	Types	Texts
		<p>1. Selection of Scalable Compact Language Models (CLMs): Identify and evaluate existing CLMs suitable for adaptation, emphasizing models with minimal computational requirements.</p> <p>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.</p> <p>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.</p> <p>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.</p> <p>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.</p> <p>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.</p> <p>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.</p> <p>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.</p>
	<b>Method</b>	<p>1. Selection and Preparation:</p> <ul style="list-style-type: none"> <li>- 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.</li> <li>- 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.</li> </ul> <p>2. Training and Adaptation Procedure:</p> <ul style="list-style-type: none"> <li>- 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.</li> <li>- 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.</li> <li>- Integrate domain-adaptive components and meta-learning algorithms into the CLMs to enable domain-specific adaptations at test time.</li> </ul> <p>3. Iterative Evaluation and Refinement:</p> <ul style="list-style-type: none"> <li>- 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.</li> <li>- Perform robustness tuning and cross-domain evaluations for each CLM and language adaptation to ensure generalizability and prevent overfitting.</li> </ul> <p>4. Human-In-The-Loop Assessment:</p> <ul style="list-style-type: none"> <li>- Enlist native speakers and domain experts to evaluate the relevance and accuracy of the model’s QA outputs for each language.</li> <li>- Incorporate feedback into the iterative training process, refining and re-adapting the models accordingly.</li> </ul> <p>5. Open-Sourcing and Community Feedback:</p> <ul style="list-style-type: none"> <li>- Make the TTT protocol, adaptive CLMs, evaluation benchmarks, and any synthetic datasets publicly available for the research community.</li> </ul> <p>6. Experiment Monitoring and Documentation:</p> <ul style="list-style-type: none"> <li>- Record all the parameters, datasets, model configurations, and evaluation metrics meticulously to ensure robustness and reproducibility.</li> <li>- Document any challenges faced, unexpected results, or adaptations made during the experiment for open-sourcing purposes.</li> </ul> <p>7. Data Analysis and Reporting:</p> <ul style="list-style-type: none"> <li>- Analyze the collected performance data quantitatively, using appropriate statistical methods to compare with non-adaptive baselines.</li> <li>- Report qualitative findings from human-in-the-loop evaluations, interpreting the implications for language model performance in low-resource language domains.</li> </ul>

*Continued on the next page*

Table 15 – Continued from the previous page

Index	Types	Texts
3	Input	<p><b>Title:</b> Whole-brain annotation and multi-connectome cell typing quantifies circuit stereotypy in <i>Drosophila</i></p> <p><b>Abstract:</b> 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 paper<sup>1</sup>, 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 database<sup>2</sup>. Crucially, this resource includes 4,552 cell types. 3,094 are rigorous consensus validations of cell types previously proposed in the "hemibrain" connectome<sup>3</sup>. 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 &gt;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><b>Entities:</b> ['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>
		<p><b>Problem</b></p> <p>Investigating the Functional Implications of Connectome Variability in <i>Drosophila</i>'s Learning and Memory Circuits Across Different Environmental and Genetic Contexts</p> <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"> <li>1. Connectome Mapping and Variability Analysis: <ol style="list-style-type: none"> <li>a. Utilize the Virtual Fly Brain database to identify and compare individual connectomes, focusing on the mushroom body.</li> <li>b. Quantify the variability in connection weights and cell type counts using statistical methods and machine learning algorithms to identify patterns of variability.</li> </ol> </li> <li>2. Behavioral Assays: <ol style="list-style-type: none"> <li>a. Design a series of learning and memory tasks for <i>Drosophila</i>, such as olfactory conditioning or visual pattern recognition.</li> <li>b. Test groups of flies with known connectome profiles under controlled environmental conditions to establish baseline behavioral data.</li> </ol> </li> <li>3. Environmental and Genetic Perturbations: <ol style="list-style-type: none"> <li>a. Expose different groups of flies to varied learning paradigms and sensory inputs to create environmental perturbations.</li> <li>b. Use CRISPR-Cas9 technology to introduce targeted mutations in genes like PINK1 or Parkin, creating genetic perturbations.</li> <li>c. Assess the impact of these perturbations on connectome structure using high-resolution imaging and reconstruction techniques.</li> </ol> </li> <li>4. Transcriptomic and Spatial Analysis: <ol style="list-style-type: none"> <li>a. Apply single-cell RNA sequencing and spatial transcriptomics to profile gene expression changes in response to environmental and genetic perturbations.</li> <li>b. Correlate transcriptomic data with connectome changes to identify molecular pathways associated with structural and functional plasticity.</li> </ol> </li> <li>5. Computational Modeling and Network Analysis: <ol style="list-style-type: none"> <li>a. Develop computational models to simulate the effects of connectome variability on neural circuit function.</li> <li>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.</li> </ol> </li> <li>6. Synthesis and Validation: <ol style="list-style-type: none"> <li>a. Integrate findings from behavioral assays, connectome mapping, transcriptomic analysis, and computational modeling.</li> <li>b. Validate the proposed models and hypotheses through iterative experimentation and refinement.</li> </ol> </li> </ol>
	Method	

*Continued on the next page*

Table 15 – Continued from the previous page

Index	Types	Texts
		<ol style="list-style-type: none"> <li>1. Pre-experimental Setup:               <ol style="list-style-type: none"> <li>a. Establish a <i>Drosophila melanogaster</i> breeding program to ensure a consistent supply of genetically similar flies for the experiments.</li> <li>b. Develop a standardized protocol for rearing flies under controlled conditions to minimize pre-experimental variability.</li> <li>c. Select and validate CRISPR-Cas9 constructs for targeted gene editing of PINK1, Parkin, and other genes of interest.</li> <li>d. Train personnel in the use of the Virtual Fly Brain database and relevant computational tools for connectome analysis.</li> </ol> </li> <li>2. Connectome Mapping and Variability Analysis:               <ol style="list-style-type: none"> <li>a. Randomly assign individual flies to either a control group or various treatment groups (environmental and genetic perturbations).</li> <li>b. Utilize high-resolution imaging techniques to map the connectomes of flies from each group, with a focus on the mushroom body.</li> <li>c. Apply statistical and machine learning algorithms to quantify and compare the variability in connection weights and cell type counts across groups.</li> </ol> </li> <li>3. Behavioral Assays:               <ol style="list-style-type: none"> <li>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.</li> <li>b. Test flies from each group in the behavioral tasks and record performance metrics.</li> <li>c. Analyze behavioral data to establish correlations with connectome profiles.</li> </ol> </li> <li>4. Environmental and Genetic Perturbations:               <ol style="list-style-type: none"> <li>a. Expose flies to different learning paradigms and sensory inputs to induce environmental perturbations.</li> <li>b. Perform gene editing using CRISPR-Cas9 to create genetic perturbations in the treatment groups.</li> <li>c. Re-map connectomes post-perturbation to assess structural changes.</li> </ol> </li> <li>5. Transcriptomic and Spatial Analysis:               <ol style="list-style-type: none"> <li>a. Collect brain tissue from flies post-behavioral assays and perform single-cell RNA sequencing and spatial transcriptomics.</li> <li>b. Analyze transcriptomic data to identify gene expression changes and correlate these with observed connectome and behavioral variations.</li> </ol> </li> <li>6. Computational Modeling and Network Analysis:               <ol style="list-style-type: none"> <li>a. Develop computational models to simulate the impact of observed connectome variability on neural circuit function.</li> <li>b. Use network analysis to integrate behavioral, connectomic, and transcriptomic data, focusing on information flow and circuit dynamics.</li> </ol> </li> <li>7. Synthesis and Validation:               <ol style="list-style-type: none"> <li>a. Integrate findings across all experimental components to formulate a cohesive understanding of the functional implications of connectome variability.</li> <li>b. Validate models and refine hypotheses through additional targeted experiments, informed by initial findings.</li> </ol> </li> </ol>

**Experiment**