

---

# Large Language Models for Automated Open-domain Scientific Hypotheses Discovery

---

Zonglin Yang<sup>1</sup> Xinya Du<sup>2</sup> Junxian Li<sup>1</sup> Jie Zheng<sup>3</sup> Soujanya Poria<sup>4</sup> Erik Cambria<sup>1</sup>

## Abstract

Hypothetical induction is recognized as the main reasoning type when scientists make observations about the world and try to propose hypotheses to explain those observations. Past research on hypothetical induction is under a constrained setting: (1) the observation annotations in the dataset are carefully manually handpicked sentences (resulting in a close-domain setting); and (2) the ground truth hypotheses are mostly commonsense knowledge, making the task less challenging. In this work, we tackle these problems by proposing the first dataset for social science academic hypotheses discovery, with the final goal to create systems that automatically generate valid, novel, and helpful scientific hypotheses, given only a pile of raw web corpus. Unlike previous settings, the new dataset requires (1) using open-domain data (raw web corpus) as observations; and (2) proposing hypotheses even new to humanity. A multi-module framework is developed for the task, including three different feedback mechanisms to boost performance, which exhibits superior performance in terms of both GPT-4 based and expert-based evaluation. To the best of our knowledge, this is the first work showing that LLMs are able to generate novel (“not existing in literature”) and valid (“reflecting reality”) scientific hypotheses<sup>1</sup>.

## 1. Introduction

Logical reasoning is central to human cognition (Goel et al., 2017). It is widely recognized as consisting of three com-

<sup>1</sup>Nanyang Technological University <sup>2</sup>University of Texas at Dallas <sup>3</sup>Huazhong University of Science and Technology <sup>4</sup>Singapore University of Technology and Design. Correspondence to: Zonglin Yang <zonglin.yang@ntu.edu.sg>.

Originally accepted by ACL 2024, to present at the ICML 2024 AI for Science workshop. Copyright by the authors.

<sup>1</sup>Dataset, code, and generated hypotheses are available at <https://github.com/ZonglinY/MOOSE.git>.

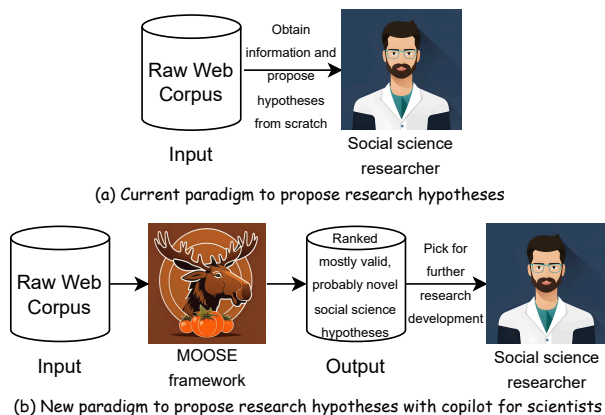


Figure 1. Comparison of the two paradigms for scientific hypotheses formulation. The new paradigm shows the role of the MOOSE framework (scientist’s copilot) and the new task setting of hypothetical induction.

ponents, which are deductive, inductive, and abductive reasoning (Yang et al., 2023b). Hypothetical induction is considered to be an important sub-type of inductive reasoning (Norton, 2003). It is recognized as the main reasoning type when scientists make observations about the world and try to propose hypotheses to explain the observations.

For example, the proposal of Geocentrism, Heliocentrism, and Newton’s law of universal gravitation based on the observations of the motion of (celestial) objects can be seen as a result of hypothetical induction. Hypothetical induction is a process of knowledge exploration from observations to hypotheses: it is challenging because it involves the exploration of knowledge that is even new to humanity. Recent research on this has two main limitations (Yang et al., 2024). Firstly, the observations in their dataset have already been manually selected from the raw web corpus, resulting in a close-domain setting. As a result, a developed system for this dataset relies on already manually selected observations, and cannot utilize the vast raw web corpus to propose hypotheses. Secondly, the ground truth hypotheses are mostly commonsense knowledge (e.g., Newton’s law), making the task less challenging since LLMs might have already seen them during pretraining. To this end, we propose a new

task setting of hypothetical induction, which is to generate novel and valid research hypotheses targeting being helpful to researchers while only given (vast) raw web corpus (Figure 1).

This hypothesis formation process is seen as the first step for scientific discovery (Wang et al., 2023a). We call this task as “auTOMated open-doMAin hypoThetical inductiOn (TOMATO)”. It is “automated” since a method for this task should automatically propose hypotheses with few human efforts; It is open-domain since it is not restricted by any manually collected data.

For the TOMATO task, we constructed a dataset consisting of 50 recent social science papers published after January 2023 in top social science journals. For each paper, social science experts collect its main hypothesis, identify its background and inspirations, find semantically similar contents for its background and inspirations from the web corpus, collect the full passage for each matched content, and use all collected web passages as raw web corpus. Although the new dataset involves many manual selection processes, the manually selected contents are used more as benchmarking human performance for comparison. In the TOMATO task, a method is required to only utilize the raw web corpus in the dataset to propose hypotheses. In addition, the raw web corpus is mostly from common news, Wikipedia, and business reviews, which means it can easily expand in scale without much human involvement.

To tackle the TOMATO task, we develop a multi-module framework called MOOSE based on large language model (LLM) prompting (Figure 4). To further improve the quality of the generated hypotheses, we also propose three different feedback mechanisms (present-feedback, past-feedback, and future-feedback) to use LLMs to retrospect and improve the LLM-generated hypotheses for better quality. For present-feedback, the intuition is that, for some modules, their generation can be evaluated by other LLMs and be provided with feedback, which can be utilized by the modules to refine their generation by taking the feedback and previous generation as input and generating again. Some modules can have feedback instantly after their generation to improve themselves. But just like the reward mechanism in reinforcement learning, some rewards (feedback) might be hard to obtain instantly, but need to wait for feedback for a future module. Similarly, we develop past-feedback where a module can benefit from the feedback for a future module. The last one is future-feedback, where a current module can provide justifications for the current module’s generation to help a future module’s generation, or can provide some initial suggestions which a future module can build upon to further provide more in-depth generation.

Both GPT-4 (OpenAI, 2023) evaluation and expert evaluation indicate that MOOSE performs better than an

LLM (Ouyang et al., 2022) based baseline, and each of the three feedback mechanisms can progressively improve the base framework. During expert evaluation, many hypotheses generated by MOOSE are recognized by social science researchers to be both novel (“not existing in the literature”) and valid (“reflecting reality”). To the best of our knowledge, this is the first work showing that LLMs can be leveraged to generate novel and valid research hypotheses, indicating the potential for LLMs to serve as a “copilot” for scientists.

## 2. Related Work

### 2.1. NLP Methods for Scientific Discovery

Zhong et al. (2023) propose a dataset where each data consists of a research goal, a corpus pair, and a discovery. However, (1) their task needs a human-provided research goal and a pairwise corpus for discovery, which is not an automated setting and has a limited application scope; (2) the discovery is not from recent publications. Wang et al. (2023b) is a concurrent work of ours. Compared the first version of two papers, they do not have an iterative feedback for novelty, reality, and clarity. Later they add for novelty, but still lack the other two. These aspects are required by inductive reasoning, and there’s an implicit trade-off between reality and novelty. Only stressing on novelty might lead to incorrect and vague generation. Bran et al. (2023) focuses on integrating computational tools in the chemistry domain, but not on providing novel chemistry findings or hypotheses. Boiko et al. (2023) focuses on using LLMs to design, plan, and execution of scientific experiments, but not on finding novel hypotheses.

### 2.2. LLM-based Self Feedback

Self-refine (Madaan et al., 2023) investigates feedback but it only focuses on present-feedback (our framework also proposes past-feedback and future-feedback), and it is not specially designed for inductive reasoning tasks. Other similar works to self-refine (Press et al., 2022; Peng et al., 2023; Yang et al., 2022; Shinn et al., 2023) also only focus on present-feedback, and their feedback is not multi-aspect nor iterative compared to ours.

Our present-feedback is developed upon a multi-aspect over-generate-then-filter mechanism (Yang et al., 2024). However, they only utilize LLMs to “filter” but not to provide feedback.

## 3. Dataset Collection

In this section, we take one publication (Gao et al., 2023) in our dataset as an example to illustrate the dataset collection process. In total, there are 50 papers published after January

**Hypothesis 2.** *Customers whose preceding customers use FR payment technology are more likely to use FR payment technology than those preceding customers do not use FR payment technology.*

Figure 2. A selected hypothesis in a social science publication collected in our dataset.

2. Hypothesis Development 2.2. Herding Effect

Figure 3. Hypothetical development section and a particular theory subsection for developing hypotheses.

2023. Table 1 shows the statistics of the subject distribution. Most social science publications highlight their hypotheses. Figure 2 shows our selected main hypothesis in the example publication. The research backgrounds are given in the introduction section. In this example paper, the background is about facial recognition payment technology’s usage in society. Most social science publications also have a “Hypothesis Development” section (some may call it by other names, e.g., “Theoretical Development”). For example, the left part (“Hypothesis Development”) in Figure 3 shows the title of this section in the example paper. In this section, several theories used to develop the main hypothesis are separately introduced. Usually, each theory takes one subsection. For example, the right part (“Herding Effect”) in Figure 3 shows the title of a subsection, which is a particular theory being used as an inspiration, which with the background can develop the hypothesis in Figure 2.

For each publication in our dataset, we identify its main hypothesis, research background, and inspirations, where the background and inspirations together provide enough information to be possible to develop the hypothesis. We also abstract the reasoning process from background and inspirations to hypothesis and note it down for each publication in our dataset. In this selected example, the reasoning process is easy, but it has medium difficulty for researchers to associate the inspiration (herding effect) to the background. For each publication, we include an expert-evaluated complexity for both the reasoning process and the association of the inspiration to the background (details in §A.3).

Instead of directly copying the background and inspirations

Social Science	Communication	5
	Psychology	7
	Human Resource Management	8
	Information System	8
	International Business	5
	Management	6
	Marketing	11

Table 1. Statistics of subject distribution of the dataset.

from the paper to construct the dataset, we try to find semantically similar text contents from the web corpus as a substitution to avoid data contamination and fit the requirement of TOMATO task that a system should propose novel and valid research hypotheses only given raw web corpus. In the example paper, we find news sentences reporting the usage of facial recognition payment as ground truth background and a Wikipedia description of the herding effect as ground truth inspiration. We also collect the web link and the full text of the manually selected web passages for backgrounds and inspirations to be used as raw web corpus.

In addition, we collect the link and the publication date for all fifty papers. We also collected fourteen survey papers in related fields that might help check the novelty of the hypotheses. The dataset is fully constructed by a social science PhD student. We illustrate why the dataset shouldn’t be collected by automatic methods in §A.4.

4. Methodology

In general, our method consists of a base multi-module framework and three feedback mechanisms (past-feedback, present-feedback, and future-feedback). We call the full framework as Multi-mOdule framewOrk with paSt present future feEdback (MOOSE). The base framework without any feedback is called MOOSE-base. MOOSE is described in Figure 4 and Algorithm 1.

4.1. Base Framework

The base framework is developed based on the intuitive understanding of how social science researchers propose an initial research hypothesis.

Firstly, a researcher needs to find a suitable research background, e.g., facial recognition payment system’s impact. This background should be proposed with a deep understanding of the societal world. Accordingly, we develop a background finder module, which reads through raw web corpus to find reasonable research backgrounds.

Secondly, since the proposed hypothesis should be novel, directly copying from raw web corpus usually is not enough. A good social science hypothesis should contain an independent variable and a dependent variable, and describe how the independent variable can influence the dependent variable. Therefore, building connections between two variables that have not been known for established connections contributes to a novel hypothesis. We hypothesize that proper inspiration can help this connection-building process, since it might serve as one of the variables itself, or might help to find such variables. However, it could consume lots of computing resources and even be practically impossible if the framework searches over the full web corpus for every found background. Nevertheless, it could be much more vi-

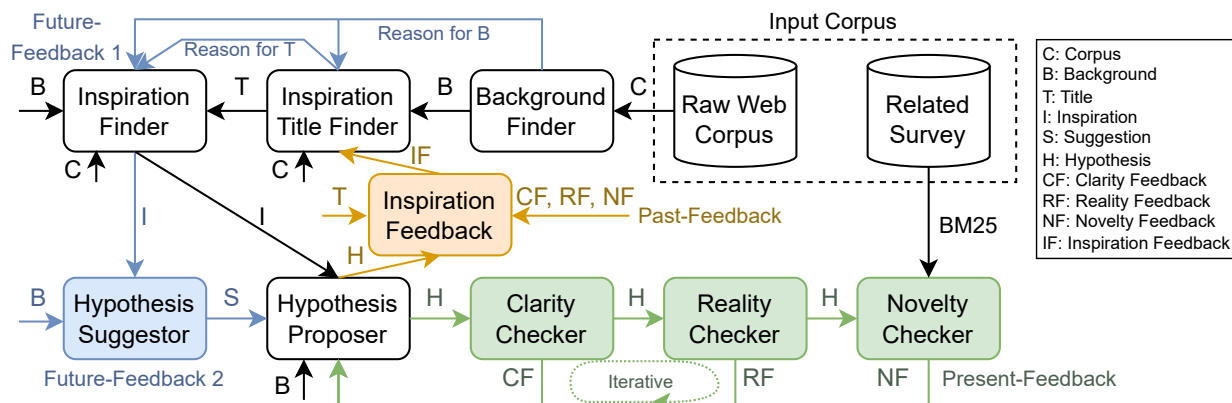


Figure 4. MOOSE: Our multi-module framework for TOMATO task. The black part is the base framework; orange part represents past-feedback.; green part represents present-feedback; blue part represents future-feedback. Each capitalized letter represents the generation of one of the modules. The same capitalized letter represents the same regardless of its color. If a module has an input arrow pointing in with a capitalized letter, it represents that this module utilizes one of its previous modules’ generation (which has the same letter pointing out) as input.

able if only searching over the titles of the corpus, and then only finding inspirational sentences in the passages which match the selected titles. Accordingly, we develop an inspiration title finder module and an inspiration finder module, together to find proper inspirations given a background.

Lastly, a hypothesis proposer module can utilize backgrounds and inspirations for hypotheses.

In general, MOOSE-base consists of a list of serializable generation modules  $M_0, M_1, \dots, M_n$  that function sequentially. The input of a module  $M_i$  is from the output of previous modules  $M_{j,j < i}$  and a raw web corpus  $C$  (and optionally a related survey corpus).  $M_i$ ’s output is represented as  $o_i$ . Feedback to  $o_i$  is represented as  $f_i$ .

#### 4.2. Present-Feedback

LLMs are not perfect and can lead to flaws in the generation, especially for those modules that undertake a difficult task. Previous work on hypothetical induction (Yang et al., 2024) tackles this problem by leveraging LLMs to identify flaws in the generation and filters those with huge flaws. Here we take a step further that instead of filtering, LLMs are leveraged to provide feedback, so that a generation can be improved rather than just filtered.

Accordingly, we define *present-feedback* as when an output  $o_i$  can be directly evaluated and provided feedback  $f_i$  (by LLMs or experts, here we use LLMs) in terms of some aspects,  $o_i$  and  $f_i$  are used as additional inputs to  $M_i$ , so that  $M_i$  can regenerate  $o_i$  to refine the previous one with  $f_i$ .

We implement present-feedback on the *Hypotheses Proposer* module, since it is a key module that undertakes a very difficult task. In terms of what aspects should the

feedback focus on, Yang et al. (2024) propose four aspects according to the philosophical definition and requirement for hypothetical induction (Norton, 2003). The aspects are whether the hypothesis (1) is consistent with observations; (2) reflects reality; (3) generalizes over the observations; (4) is clear and meaningful.

In MOOSE, we basically adopt the four aspects but reframe them to better fit the current task. Specifically, aspect (2) contains aspect (1) most of the time (unless the observations are wrongly described). To save computing power, we adopt aspect (2) but not aspect (1). In addition, we reframe aspect (3) as whether the hypothesis is novel, and reframe aspect (4) as whether the hypothesis is clear and provides enough details. Accordingly, we develop a reality checker module, a novelty checker module, and a clarity checker module in Figure 4.

#### 4.3. Past-Feedback

Just like the reward mechanism in reinforcement learning, some modules’ generation can only be evaluated at a future time point. For instance, it is hard to give feedback on the selected inspirations unless we know what hypotheses these inspirations could lead to. Accordingly, we develop *past-feedback* as when it is hard to directly evaluate  $o_i$ , the framework continues to run until generating  $o_{j,j > i}$ , where  $o_j$  is highly influenced by  $o_i$  and can be directly evaluated to obtain present-feedback  $f_j$ . Then  $o_i, o_j$ , and  $f_j$  are utilized, possibly by an additional module implemented with an LLM, to provide past-feedback  $f_i$  to  $M_i$ , so that  $M_i$  can regenerate  $o_i$  with  $f_i$  to refine the previous  $o_i$ .

We implement past-feedback on the *Inspiration Title Finder* module. The intuition is that improper inspirations can lead

to low-quality hypotheses, and it is hard to directly evaluate inspirations.

#### 4.4. Future-Feedback

We also develop *future-feedback*, targeting at providing additional useful information for a future module  $M_j$  to generate  $o_j$  in better quality. Specifically, we develop future-feedback-1 (FF1) and future-feedback-2 (FF2). FF1 is that in addition to  $o_i$ , justifications (reasons) of  $o_i$  are also provided to  $M_{j,j>i}$  so that  $M_j$  can better leverage  $o_i$ ; FF2 is that for a key module  $M_j$  that handles a very complex task, an additional module  $M_{j-0.5}$  is being placed before  $M_j$ , so that  $M_{j-0.5}$  can undertake some of the reasoning burdens of  $M_j$  to improve the quality of  $o_j$ . For example, in MOOSE,  $M_{j-0.5}$  is to provide preliminary suggestions for  $M_j$ .

Specifically in the MOOSE framework, for FF1, no additional modules are needed. Instead, we modify the prompt to require  $M_i$  to not only generate  $o_i$  but also provide the justification of  $o_i$ . We implement it on the *Background Finder* and the *Inspiration Title Finder* modules. The intuition is that it could be helpful if the *Inspiration Title Finder* module knows not only the background but also what possible research topics could be conducted for this background so as to select suitable titles; it could be also helpful for the *Inspiration Finder* module to know why this background was selected and what potentially helpful inspirations could be found from the passage with the corresponding selected titles. For FF2, we implement it on the *Hypothesis Proposer* module, since proposing hypotheses is a very important and complex task. Accordingly, we develop a *Hypothesis Suggestor* module (as  $M_{j-0.5}$ ) to provide some initial suggestions on how to utilize the inspirations and background first, and then *Hypothesis Proposer* (as  $M_j$ ) can build upon the suggestions to generate more novel and more complicated hypotheses.

## 5. Experiments

### 5.1. Evaluation Metrics & Details

We conduct both automatic evaluation and human evaluation for the experiments.

For automatic evaluation, we adopt validness, novelty, and helpfulness as three aspects for GPT-4 to evaluate. We choose validness and novelty because they are the two basic requirements for hypothetical induction illustrated in philosophical literature (Norton, 2003; Yang et al., 2024). In addition, these two scores also highly resemble the current ACL review form, which requires reviewers to score submitted papers on soundness and excitement aspects. We choose helpfulness because the final goal of the TOMATO task is to provide help and assistance for human scientists.

	Validness	Novelty	Helpfulness
Baseline	3.954	2.483	3.489
MOOSE-base	3.907	3.081	3.859
w/ <i>future-feedback</i>	<b>3.955</b>	3.226	<b>3.953</b>
w/ <i>future-</i> and <i>past-feedback</i>	3.916	<b>3.390</b>	3.931

Table 2. Effect of MOOSE-base, *future-feedback* and *past-feedback* (evaluated by GPT-4). MOOSE-related results are averaged over iterations of *present-feedback*. Base model is GPT-3.5.

	Validness	Novelty	Helpfulness
MOOSE (w/o <i>present-feedback</i> )	3.823	3.114	3.809
w/ 1 iteration of <i>present-feedback</i>	3.918	3.199	3.900
w/ 2 iterations of <i>present-feedback</i>	3.951	3.293	3.956
w/ 3 iterations of <i>present-feedback</i>	3.969	3.270	<b>3.962</b>
w/ 4 iterations of <i>present-feedback</i>	<b>3.970</b>	<b>3.329</b>	3.951

Table 3. Effect of *present-feedback* (evaluated by GPT-4). Base model is GPT-3.5.

In §A.5 we illustrate why we don’t adopt evaluation metrics such as (1) relevance and significance, and (2) BLEU (Papineni et al., 2002), ROUGE (Lin, 2004), or METEOR (Banerjee & Lavie, 2005).

For human (expert) evaluation, evaluation metrics are the same. Three experts (social science PhD students) take charge of the expert evaluation. They evaluate on 400 randomly selected hypotheses from the baseline and variants of the MOOSE framework. To avoid any bias, they are not told which methods we are comparing; the order of generated hypotheses to compare is also randomized. We introduce how the 400 hypotheses are selected in §A.6, and the high expert agreement in §A.7.

Each metric is on a 5-point Likert scale. Both experts and GPT-4 are given the same description of the scale and evaluation standard of the three aspects (listed in §A.9).

Out of the metrics, we consider the novelty metric to be relatively more important than the validness metric. Because the goal of the TOMATO task is to assist human researchers, but not to directly add the machine-proposed hypotheses to the literature. If the hypotheses are fully valid but not novel, then they are not helpful at all; but if the hypotheses are novel but not valid, then they can still be possible to inspire human researchers to develop novel and valid hypotheses. Helpfulness is also an important metric since it could be seen as an overall evaluation of a hypothesis.

In §A.8, we introduce the surprisingly high consistency between expert evaluation and GPT4 evaluation, indicating that GPT-4 might be able to provide a relatively reliable evaluation for machine-generated social science hypotheses.

	Validness	Novelty	Helpfulness
Baseline	3.579	2.276	2.632
MOOSE-base	3.500	2.855	3.026
w/ <i>future-feedback</i>	3.645	3.105	3.303
w/ <i>future- and past-feedback</i>	<b>3.750</b>	<b>3.197</b>	<b>3.368</b>

Table 4. Effect of MOOSE-base, *future-feedback* and *past-feedback* (evaluated by *experts*). MOOSE results are selected from the 5<sup>th</sup> iteration of *present-feedback*. Base model is GPT-3.5.

	Validness	Novelty	Helpfulness
MOOSE-base (w/o <i>present-feedback</i> )	3.342	2.382	2.500
w/ 2 iterations of <i>present-feedback</i>	<b>3.539</b>	2.803	2.934
w/ 4 iterations of <i>present-feedback</i>	3.500	<b>2.855</b>	<b>3.026</b>
MOOSE (w/o <i>present-feedback</i> )	3.224	2.737	2.855
w/ 2 iterations of <i>present-feedback</i>	3.579	<b>3.250</b>	3.342
w/ 4 iterations of <i>present-feedback</i>	<b>3.750</b>	3.197	<b>3.368</b>

Table 5. Effect of *present-feedback* (evaluated by *experts*). Base model is GPT-3.5.

## 5.2. Baselines & Base Model Selection

Since the TOMATO task is to propose hypotheses given only corpus, a natural baseline is to use a corpus chunk as input, and directly output hypotheses.

Except for §6.3, we use `gpt-3.5-turbo` for each module in MOOSE. To be fair, the baseline is also instantiated with `gpt-3.5-turbo`. The training data of the model checkpoint is up to September 2021, while all papers in our dataset are published after January 2023, so the model has not seen any of the collected papers in the dataset. In §6.3, we investigate the effect of base model selection by using Claude3-Opus (Anthropic, 2024) for each module in MOOSE.

## 5.3. Main Results

In this subsection, we compare MOOSE-base with the baseline and examine the effect of each of the three feedback mechanisms to MOOSE-base.

We first introduce the number of generated hypotheses being evaluated in §5.3 and §6. For experiments evaluated with GPT-4, fifty backgrounds are selected for each method. For MOOSE-related methods, for each background, on average around 6 inspirations are extracted, resulting in 4 different hypotheses. Each hypothesis leads to another 4 more refined ones with *present-feedback*. Therefore on average for each MOOSE-related method in GPT-4 evaluation tables, around  $50 \times 4 \times 5 = 1000$  hypotheses are evaluated. For experiments evaluated with expert evaluation, in general, we randomly select one hypothesis for each background, resulting in 50 hypotheses evaluated for each line of the method in expert evaluation tables.

Table 2 shows GPT-4’s evaluation targeting at comparing MOOSE-base and the baseline and shows the effect of *future-feedback* and *past-feedback*. In this table, MOOSE-related results are averaged over iterations of *present-feedback* to not be influenced by *present-feedback*. MOOSE-base largely outperforms the baseline in terms of both novelty and helpfulness, but slightly lower in terms of validness. As illustrated in §5.1, since the purpose of the TOMATO task is to inspire and help human researchers, novelty and helpfulness metrics should be more important. In practice, we find many hypotheses from baseline almost only rephrasing some sentences in the input corpus, adding little novelty content. MOOSE-base with *future-feedback* comprehensively outperforms MOOSE-base in terms of all three metrics. MOOSE-base with both *future* and *past-feedback* largely outperforms MOOSE-base with *future-feedback* in novelty and performs slightly lower in validness and helpfulness metrics. One of the reasons is that the *past-feedback* may focus more on the novelty aspect because the novelty checker module provides more negative *present-feedback* than the reality checker module.

Table 3 shows the effect of *present-feedback* with GPT-4 evaluation. In this table, the results are averaged over three experiments: MOOSE-base, MOOSE-base with *future-feedback*, and MOOSE-base with both *future* and *past-feedback* to focus on *present-feedback*. It shows that as more iterations of *present-feedback* are conducted, validness and novelty steadily go up; helpfulness also steadily goes up but reaches the best performance with 3 iterations of *present-feedback*.

Table 4 shows expert evaluation results on the comparison between MOOSE-base and the baseline, and the effect of *future-feedback* and *past-feedback*. MOOSE-related results are selected from the 5<sup>th</sup> iteration of *present-feedback*. Similar to GPT-4 evaluation, MOOSE-base largely outperforms the baseline in terms of Novelty and Helpfulness; MOOSE-base with *future-feedback* comprehensively outperforms MOOSE-base. Different from GPT-4 evaluation, MOOSE-base with *future* and *past-feedback* also comprehensively outperforms MOOSE-base with *future-feedback*. We think one of the reasons could be that GPT-4 might grade validness based on how frequently it has seen relevant texts, but not true understanding of the world. Therefore a more novel hypothesis might tend to have a relatively lower score in validness and helpfulness under GPT-4 evaluation.

Table 5 shows the expert evaluation of *present-feedback*. MOOSE-base and MOOSE are both evaluated. Overall performance generally goes up with more iterations of *present-feedback*, but there might be an optimal number of iterations.

	Validness	Novelty	Helpfulness
Rand background	3.954	2.483	3.489
Rand background and rand inspirations	3.773	2.957	3.643
Rand background and BM25 inspirations	3.585	<b>3.364</b>	3.670
GPT-3.5 picked background and inspirations	<b>3.812</b>	2.818	<b>3.733</b>
Groundtruth background and inspirations	<b>3.876</b>	3.000	3.806
Groundtruth hypotheses	3.700	<b>3.380</b>	<b>3.880</b>

Table 6. Analysis of retrieval’s effect on generated hypotheses (evaluated by GPT-4). No methods here utilize any feedback mechanisms. Base model is GPT-3.5.

## 6. Analysis

### 6.1. Background and Inspirations

Here we try to answer “Is ChatGPT necessary for background and inspiration selection?”.

Table 6 shows various methods for background and inspiration selection. In general, there might be a validness-novelty trade-off that if a method reaches a high novelty score, then it is usually hard for it to reach a high validness score. It is surprising that a randomly selected background and randomly selected inspirations can lead to hypotheses with relatively comparable validness and novelty to ChatGPT-picked background and inspirations. Empirically we hypothesize the reason is that randomly picked inspirations are mostly not related to the background, resulting in a high novelty (but less validness and helpfulness). In addition, BM25 (Robertson et al., 2009) picked background and inspirations reach a much higher novelty score compared to ChatGPT-picked ones. Empirically we do not find BM25 retrieved inspirations to be similar to the background, but they are usually with more concrete contents compared with random inspirations. Not surprisingly, ChatGPT picked background and inspirations reach the highest helpfulness score among those without any ground-truth annotations. Lastly, ground-truth hypotheses reach the highest novelty and helpfulness.

### 6.2. More Ablation Studies

Table 7 shows ablation studies on future-feedback, access to surveys, and the selection of corpus.

Firstly, for future-feedback, we separately test the effect of FF1 and FF2. Without FF2, performance comprehensively drops; without FF1, performance drops on validness and novelty, with helpfulness remaining comparable. It seems that FF2 is more significant than FF1. However, the fact that FF1 works on inspiration title finder and inspiration finer modules does not mean that it works on all modules. Empirically we find that adding the reasons (or prospects) for background and inspirations to the hypothesis proposer module will cause a more valid but much less novel generation of hypotheses. The reason is that the hypothesis proposer

	Validness	Novelty	Helpfulness
MOOSE	3.916	3.390	3.931
w/o future-feedback-2	3.895	3.281	3.918
w/o future-feedback-1	3.882	3.355	3.935
w/o access to related survey	3.889	<b>3.431</b>	3.886
w/ randomized corpus	<b>3.941</b>	3.227	<b>3.955</b>

Table 7. More ablation study (evaluated by GPT-4). Results are averaged over iterations of *present-feedback*. Base model is GPT-3.5.

module tends to simply follow the prospects, which do not have a global view of both background and all inspirations, but only focus on one background or one inspiration. Instead, FF2 (the hypothesis suggestor module) has the global view and only provides soft initial suggestions on how to combine the background and inspirations together. With the hypotheses suggestor module, the hypotheses proposer module is prompted to further combine the initial suggestions and other inspirations to propose hypotheses. To be fair, MOOSE-base, which is not equipped with the hypothesis suggestor module, has the same prompt to combine the inspirations together (just without suggestions) to propose hypotheses.

Secondly, we cut the access of novelty detector to related surveys to check the effect of related surveys. As a result, novelty largely goes up (0.04), and validness goes down to around 0.26. Empirically one of the main reasons is that BM25 hardly retrieves enough similar survey chunks, so that access to the survey leads novelty detector to tend to reply the hypotheses are novel since it is not mentioned in the related survey. Without present-feedback, MOOSE and MOOSE w/o access to survey perform quite comparably.

Lastly, the raw corpus in the dataset is from two sources: passages that contain the ground truth backgrounds and passages that contain the ground truth inspirations. In all of the previous experiments, backgrounds are extracted from the background passages, and inspirations are extracted from the inspirations passages. To see whether the passages are only restricted to their designed role, in MOOSE w/ randomized corpus experiment, we use inspiration corpus for background extraction and use both inspiration and background corpus for inspiration extraction. As a result, validness goes up by about 0.025, while novelty goes down by about 0.16. We think one of the reasons is that, in this setting, after selecting a background from an inspiration passage, MOOSE tends to retrieve the same inspiration passage to find inspirations, which leads to less novel results.

### 6.3. Effect of Base Model Selection

In all previous experiments, we adopt GPT-3.5 as the base model. In this section, we investigate the effect of base model selection by using Claude3-Opus as the base

	Validness	Novelty	Helpfulness
Baseline	3.884	2.925	3.856
MOOSE-base	<b>3.967</b>	3.392	3.939
w/ future-feedback	3.926	3.694	<b>3.966</b>
w/ future- and past-feedback	3.875	<b>4.177</b>	3.868

Table 8. Effect of MOOSE-base, future-feedback and past-feedback (evaluated by GPT-4). MOOSE-related results are averaged over iterations of present-feedback. Base model is Claude3-Opus.

	Validness	Novelty	Helpfulness
MOOSE (w/o present-feedback)	3.793	3.683	3.870
w/ 1 iteration of present-feedback	3.896	3.804	3.937
w/ 2 iterations of present-feedback	3.961	3.730	3.939
w/ 3 iterations of present-feedback	<b>3.983</b>	<b>3.809</b>	<b>3.946</b>
w/ 4 iterations of present-feedback	3.980	3.757	3.930

Table 9. Effect of present-feedback (evaluated by GPT-4). Base model is Claude3-Opus.

model for each module in MOOSE.

With Claude3-Opus as the base model, we again analyze the effect of MOOSE-base, past-feedback, and future-feedback in Table 8; and analyze the effect of present-feedback in Table 9. The experiment settings of Table 8 and Table 9 are exactly the same as in Table 2 and Table 3 correspondingly, but only differ in the base model selection.

In general, there are two conclusions. Firstly, MOOSE’s components stay effective regardless of different base model selection. It shows the robustness of the MOOSE framework in terms of different base model. Secondly, the absolute evaluation scores on all three metrics largely improved with Claude3-Opus compared to GPT-3.5, indicating the even larger potential of the MOOSE framework when more powerful LLMs are available.

#### 6.4. Qualitative Analysis

The following box shows one generated counter-intuitive hypothesis (expert evaluation appended).

*In collectivist cultures, individuals engage in more conspicuous consumption behaviors compared to individualistic cultures. (Validness: 3.3; Novelty: 4.0; Helpfulness: 4.0)*

Here is the assessment from one of the experts:

*The main reason I give a high mark for both three dimensions of this hypothesis is because:*

*(1) For validness, this hypothesis is based on existing cultural theories and empirical evidence that suggests cultural values significantly impact consumer behavior. It aligns with established concepts like collectivism and individual-*

*ism that have been widely studied in cross-cultural psychology.*

*(2) For novelty, this hypothesis is counter-intuitive to some extent. Prior research has shown that collectivist cultures often prioritize group harmony, cooperation, and social cohesion over individual desires. This emphasis on collective well-being might suggest a reduced inclination toward overt displays of personal wealth or status through conspicuous consumption. However, this hypothesis suggests the opposite that collectivist culture’s members engage in more conspicuous consumption, which is more commonly linked to individualistic societies in popular perceptions. This challenges the notion that members of collectivist cultures avoid conspicuous consumption behaviors.*

*(3) For helpfulness, if this hypothesis is confirmed, it could have significant practical implications. Understanding the impact of cultural values on conspicuous consumption can assist businesses and marketers in crafting more effective cross-cultural marketing strategies. It could also aid policymakers in addressing societal issues related to consumerism.*

In addition to the analysis of this counter-intuitive example, we also provide qualitative analysis on the difference between hypotheses generated from the baseline, MOOSE-base, MOOSE-base w/ future-feedback, and MOOSE-base w/ future and past-feedback in §A.11. More qualitative analysis on highly scored generated hypotheses can be found in §A.12. Additionally, §A.13 illustrates factors for good hypotheses in social science (particularly in Business). §A.14 shows how MOOSE formulates a hypothesis by giving the generation of each of the modules in MOOSE.

## 7. Conclusion

In this paper, we propose a new task, automated open-domain hypothetical induction (TOMATO), which is the first task in NLP to focus on social science research hypotheses discovery. Along with the task, we construct a dataset consisting of 50 recent social science papers published in top academic journals. We also developed a multi-module framework MOOSE for the TOMATO task, which contains a base framework and three novel feedback mechanisms. Experiments indicate that MOOSE-base outperforms an LLM-based baseline, and the three feedback mechanisms can progressively further improve over MOOSE-base. Surprisingly, evaluated by PhD students, MOOSE is able to produce many novel (“not existing in the literature”) and valid (“reflecting reality”) research hypotheses. To the best of our knowledge, this is the first work showing that LLMs can be leveraged to generate novel and valid scientific hypotheses, indicating the potential of LLMs to serve as a “copilot” for scientists.



## References

- Anthropic, A. The claude 3 model family: Opus, sonnet, haiku. *Claude-3 Model Card*, 2024.
- Banerjee, S. and Lavie, A. METEOR: An automatic metric for MT evaluation with improved correlation with human judgments. In *Proceedings of the ACL Workshop on Intrinsic and Extrinsic Evaluation Measures for Machine Translation and/or Summarization*, pp. 65–72, Ann Arbor, Michigan, June 2005. Association for Computational Linguistics. URL <https://aclanthology.org/W05-0909>.
- Boiko, D. A., MacKnight, R., and Gomes, G. Emergent autonomous scientific research capabilities of large language models. *CoRR*, abs/2304.05332, 2023. doi: 10.48550/arXiv.2304.05332. URL <https://doi.org/10.48550/arXiv.2304.05332>.
- Bosselut, A., Rashkin, H., Sap, M., Malaviya, C., Celikyilmaz, A., and Choi, Y. COMET: Commonsense transformers for automatic knowledge graph construction. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pp. 4762–4779, Florence, Italy, July 2019. Association for Computational Linguistics. doi: 10.18653/v1/P19-1470. URL <https://aclanthology.org/P19-1470>.
- Bran, A. M., Cox, S., White, A. D., and Schwaller, P. Chemcrow: Augmenting large-language models with chemistry tools. *arXiv preprint arXiv:2304.05376*, 2023.
- Das, R., Zaheer, M., Thai, D., Godbole, A., Perez, E., Lee, J. Y., Tan, L., Polymenakos, L., and McCallum, A. Case-based reasoning for natural language queries over knowledge bases. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pp. 9594–9611, Online and Punta Cana, Dominican Republic, November 2021. Association for Computational Linguistics. doi: 10.18653/v1/2021.emnlp-main.755. URL <https://aclanthology.org/2021.emnlp-main.755>.
- Gao, J., Rong, Y., Tian, X., and Yao, Y. Improving convenience or saving face? an empirical analysis of the use of facial recognition payment technology in retail. *Information Systems Research*, 2023.
- Goel, V., Navarrete, G., Noveck, I. A., and Prado, J. The reasoning brain: The interplay between cognitive neuroscience and theories of reasoning, 2017.
- Lin, C.-Y. Rouge: A package for automatic evaluation of summaries. In *Text summarization branches out*, pp. 74–81, 2004.
- Liu, C.-W., Lowe, R., Serban, I., Noseworthy, M., Charlin, L., and Pineau, J. How NOT to evaluate your dialogue system: An empirical study of unsupervised evaluation metrics for dialogue response generation. In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, pp. 2122–2132, Austin, Texas, November 2016. Association for Computational Linguistics. doi: 10.18653/v1/D16-1230. URL <https://aclanthology.org/D16-1230>.
- Madaan, A., Tandon, N., Gupta, P., Hallinan, S., Gao, L., Wiegrefe, S., Alon, U., Dziri, N., Prabhumoye, S., Yang, Y., Welleck, S., Majumder, B. P., Gupta, S., Yazdanbakhsh, A., and Clark, P. Self-refine: Iterative refinement with self-feedback. *CoRR*, abs/2303.17651, 2023. doi: 10.48550/arXiv.2303.17651. URL <https://doi.org/10.48550/arXiv.2303.17651>.
- Norton, J. D. A little survey of induction. 2003.
- OpenAI. GPT-4 technical report. *CoRR*, abs/2303.08774, 2023. doi: 10.48550/arXiv.2303.08774. URL <https://doi.org/10.48550/arXiv.2303.08774>.
- Ouyang, L., Wu, J., Jiang, X., Almeida, D., Wainwright, C. L., Mishkin, P., Zhang, C., Agarwal, S., Slama, K., Ray, A., Schulman, J., Hilton, J., Kelton, F., Miller, L., Simens, M., Askell, A., Welinder, P., Christiano, P. F., Leike, J., and Lowe, R. Training language models to follow instructions with human feedback. In *NeurIPS*, 2022. URL [http://papers.nips.cc/paper\\_files/paper/2022/hash/blefde53be364a73914f58805a001731-Abstract-Conference.html](http://papers.nips.cc/paper_files/paper/2022/hash/blefde53be364a73914f58805a001731-Abstract-Conference.html).
- Pan, F., Mulkar-Mehta, R., and Hobbs, J. R. Annotating and learning event durations in text. *Comput. Linguistics*, 37(4):727–752, 2011. doi: 10.1162/COLI\_a\_00075. URL [https://doi.org/10.1162/COLI\\_a\\_00075](https://doi.org/10.1162/COLI_a_00075).
- Papineni, K., Roukos, S., Ward, T., and Zhu, W.-J. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the 40th annual meeting of the Association for Computational Linguistics*, pp. 311–318, 2002.
- Peng, B., Galley, M., He, P., Cheng, H., Xie, Y., Hu, Y., Huang, Q., Liden, L., Yu, Z., Chen, W., and Gao, J. Check your facts and try again: Improving large language models with external knowledge and automated feedback. *CoRR*, abs/2302.12813, 2023. doi: 10.48550/arXiv.2302.12813. URL <https://doi.org/10.48550/arXiv.2302.12813>.
- Press, O., Zhang, M., Min, S., Schmidt, L., Smith, N. A., and Lewis, M. Measuring and narrowing the compositionality gap in language models. *CoRR*, abs/2210.03350,

2022. doi: 10.48550/arXiv.2210.03350. URL <https://doi.org/10.48550/arXiv.2210.03350>.
- Robertson, S., Zaragoza, H., et al. The probabilistic relevance framework: Bm25 and beyond. *Foundations and Trends® in Information Retrieval*, 3(4):333–389, 2009.
- Shinn, N., Labash, B., and Gopinath, A. Reflexion: an autonomous agent with dynamic memory and self-reflection. *CoRR*, abs/2303.11366, 2023. doi: 10.48550/arXiv.2303.11366. URL <https://doi.org/10.48550/arXiv.2303.11366>.
- Wang, H., Fu, T., Du, Y., Gao, W., Huang, K., Liu, Z., Chandak, P., Liu, S., Van Katwyk, P., Deac, A., et al. Scientific discovery in the age of artificial intelligence. *Nature*, 620(7972):47–60, 2023a.
- Wang, Q., Downey, D., Ji, H., and Hope, T. Learning to generate novel scientific directions with contextualized literature-based discovery. *CoRR*, abs/2305.14259, 2023b. doi: 10.48550/arXiv.2305.14259. URL <https://doi.org/10.48550/arXiv.2305.14259>.
- Yang, K., Tian, Y., Peng, N., and Klein, D. Re3: Generating longer stories with recursive reprompting and revision. In Goldberg, Y., Kozareva, Z., and Zhang, Y. (eds.), *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing, EMNLP 2022, Abu Dhabi, United Arab Emirates, December 7-11, 2022*, pp. 4393–4479. Association for Computational Linguistics, 2022. doi: 10.18653/v1/2022.emnlp-main.296. URL <https://doi.org/10.18653/v1/2022.emnlp-main.296>.
- Yang, Z., Du, X., Rush, A., and Cardie, C. Improving event duration prediction via time-aware pre-training. In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pp. 3370–3378, Online, November 2020. Association for Computational Linguistics. doi: 10.18653/v1/2020.findings-emnlp.302. URL <https://aclanthology.org/2020.findings-emnlp.302>.
- Yang, Z., Du, X., Cambria, E., and Cardie, C. End-to-end case-based reasoning for commonsense knowledge base completion. In *Proceedings of the 17th Conference of the European Chapter of the Association for Computational Linguistics*, pp. 3509–3522, Dubrovnik, Croatia, May 2023a. Association for Computational Linguistics. URL <https://aclanthology.org/2023.eacl-main.255>.
- Yang, Z., Du, X., Mao, R., Ni, J., and Cambria, E. Logical reasoning over natural language as knowledge representation: A survey. *CoRR*, abs/2303.12023, 2023b. doi: 10.48550/arXiv.2303.12023. URL <https://doi.org/10.48550/arXiv.2303.12023>.
- Yang, Z., Dong, L., Du, X., Cheng, H., Cambria, E., Liu, X., Gao, J., and Wei, F. Language models as inductive reasoners. In Graham, Y. and Purver, M. (eds.), *Proceedings of the 18th Conference of the European Chapter of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 209–225, St. Julian’s, Malta, March 2024. Association for Computational Linguistics. URL <https://aclanthology.org/2024.eacl-long.13>.
- Zhong, R., Zhang, P., Li, S., Ahn, J., Klein, D., and Steinhart, J. Goal driven discovery of distributional differences via language descriptions. *CoRR*, abs/2302.14233, 2023. doi: 10.48550/arXiv.2302.14233. URL <https://doi.org/10.48550/arXiv.2302.14233>.

## A. Appendix

### A.1. Hyper-parameters

Experiments in §6.3 adopts Claude3-Opus, all other experiments are conducted with gpt-3.5-turbo.

Both Claude3-Opus and gpt-3.5-turbo use 0.9 temperature and 0.9 top-p.

The hyperparameters for GPT-4 evaluation are 0.0 temperature to ensure the evaluation scores are stable, and 0.9 top-p.

### A.2. More Related Works on Reasoning and Scientific Discovery

This paper is a successive work in inductive reasoning and is different from commonsense reasoning (Bosselut et al., 2019; Yang et al., 2020) in that the novel social science hypotheses do not belong to commonsense.

Case-based reasoning (Das et al., 2021; Yang et al., 2023a) also falls in the domain of inductive reasoning, but case-based reasoning is more about high-level guidance on methodology design (case retrieve, reuse, revise, and retain), which is not involved in this paper.

### A.3. Dataset Complexity Distribution

Table 10 illustrates the complexity distribution of the proposed dataset from both reasoning and association perspectives. “Easy” in the table means it is relatively easy compared to other publications in the dataset, but does not mean it is actually easy to induce the hypotheses.

### A.4. Why the Tomato Dataset Shouldn’t Be Collected by Automatic Methods

Firstly, there are many hypotheses in a social science publication, which might need an expert to identify which hypothesis is suitable for this task (e.g., whether it is a main hypothesis, whether the background and inspirations are properly introduced).

Secondly, the background and inspirations scatter in a publication. It needs a deep domain understanding of the hypothesis, related background, and inspirations to select the background and inspirations out to form a complete reasoning chain to conclude the hypothesis.

Thirdly, it needs enough domain knowledge to find semantically similar texts (similar to the groundtruth selected background and inspirations) from the web, where the texts should contain enough details to help elicit the hypothesis.

### A.5. Why Not Using Other Evaluation Metrics

Other relevant aspects from related literature include relevance (Wang et al., 2023b) and significance (Zhong et al., 2023).

We do not adopt relevance because our task setting is the automated and open domain, without a manually given background; neither for significance because social science is different from engineering subjects — (1) every hypothesis is to reflect the reality of the world, and as long as it reflects the world, it is significant. Therefore it is hard to tell which one is more significant even by experts; (2) the evaluation standard of significance varies from time to time. For example, in the 60s, conducting research on how to improve the assembly line’s efficiency as much as possible was seen as very significant. However, in recent decades, how to alleviate the psychological depression of assembly line workers is seen as more significant.

We do not adopt BLEU (Papineni et al., 2002), ROUGE (Lin, 2004), or METEOR (Banerjee & Lavie, 2005) as evaluation metric to compare the proposed hypothesis and the ground truth hypothesis since (1) proposing novel research hypotheses is an open problem, and (2) TOMATO has an automated open domain setting, which means the automatically selected

	Reasoning Complexity	Association Complexity
Easy	24	12
Medium	17	25
Hard	9	13

Table 10. Statistics of the complexity of the dataset.

	Validness	Novelty	Helpfulness
Hard Consistency	0.298	0.337	0.361
Soft Consistency	0.755	0.793	0.791

Table 11. Hard and soft consistency scores between evaluation from different experts in terms of Validness, Novelty, and Helpfulness metrics.

	Validness	Novelty	Helpfulness
Hard Consistency	0.485	0.392	0.321
Soft Consistency	0.850	0.823	0.773

Table 12. Hard and soft consistency scores between expert evaluation and GPT-4 evaluation in terms of Validness, Novelty, and Helpfulness metrics.

background and inspirations are hardly the same as a few given ground truth ones (if background and inspirations are not the same, then it is meaningless to compare the hypothesis). Liu et al. (2016) have conducted a comprehensive analysis that they also reached a similar conclusion that BLEU, METEOR, or ROUGE is not suitable for an open-ended task (such as a dialogue system).

### A.6. Hypotheses Selection for Expert Evaluation

In total, we randomly selected 400 hypotheses to be evaluated by experts. Specifically, for each background passage in the dataset (out of 50), we use 4 methods (which are to be compared) to collect in total 8 hypotheses.

The 8 hypotheses are from (1) the baseline; (2) the MOOSE-base framework; (3) MOOSE-base + future-feedback; (4) MOOSE-base + future-feedback + past-feedback. For (2) and (4), we collect three hypotheses, which are (a) without present-feedback; (b) after 2 iterations of present-feedback; and (c) after 4 iterations of present-feedback. For (1) and (3), we only collect one hypothesis, which is without present-feedback.

With these collections, we can evaluate the effect of both the MOOSE-base framework and the three feedback methods, leading to results in Table 4 and Table 5.

Out of the three experts, one expert evaluates the full 400 hypotheses, and the other two each evaluate 104 hypotheses (the first and second 104 hypotheses out of 400). The reason we choose the number “104” is that (1) social science PhD students are quite busy and two of them can only have time to evaluate around 100 hypotheses; (2) the number should be dividable by 8 (since every 8 hypotheses form a group for comparison).

The results of the expert evaluation are averaged over the three experts. Specifically, expert evaluation essentially compares the 8 hypotheses within a group. The 400, 104, and 104 hypotheses evaluation scores can be written as arrays of [50, 8], [13, 8], and [13, 8]. We concatenate them to [76, 8], and average them across the first dimension.

The payment for expert evaluation is \$1 per hypothesis.

### A.7. Expert Qualification and Expert Agreement

The constructed dataset covers many subjects, but every collected publication is somewhat related to Marketing, which is a big topic in Business research. It is common in social science to conduct research that connects with other social science domains. The experts for expert evaluation are three PhD students majoring in Marketing. Therefore the experts are qualified enough to provide assessment for machine-generated hypotheses in the domain.

The consistency scores between experts are shown in Table 11. The soft consistency and hard consistency are defined in §A.8. All soft consistency scores are above 0.75 means, and the average difference between experts in terms of each metric is less than 1 (out of a 5-point scale), exhibiting high expert evaluation agreement.

Aspect 1: Validness	
5 points	The hypothesis completely reflects the reality.
4 points	The hypothesis almost completely reflects the reality, but has only one or two minor conflictions that can be easily modified.
3 points	The hypothesis has at least one moderate conflict or several minor conflicts.
2 points	The hypothesis has at least one major confliction with the reality or only establishes in very rare circumstances that are not mentioned in this hypothesis.
1 point	The hypothesis completely violates the reality.

Table 13. Evaluation standard for *Validness*.

Aspect 1: Novelty	
5 points	The hypothesis is completely novel and has not been proposed by any existing literature.
4 points	The main argument or several sub-arguments of the hypothesis are novel.
3 points	The main argument is not novel, only one or two sub-arguments appear to be novel.
2 points	The full hypothesis is not novel, but the way it combines the topics can be inspiring for human researchers.
1 point	The hypothesis is not novel at all and not inspiring for human researchers.

Table 14. Evaluation standard for *Novelty*.

### A.8. Consistency Between Expert Evaluation and GPT-4 Evaluation

To check the consistency between expert evaluation and GPT-4 evaluation, we use the expert evaluation results and find the corresponding GPT-4 evaluation results. In total, there are 400 hypotheses evaluated by experts, so the sample we use to calculate the consistency score is 400.

Specifically, similar to Pan et al. (2011), for soft consistency, if the absolute difference between expert evaluation and GPT-4 evaluation (both are on a 5-point scale) is 0/1/2/3/4, then we assign a consistency score of 1.00/0.75/0.50/0.25/0.00; for hard consistency, if only the difference is 0, can the consistency score be 1.00, otherwise consistency score is 0.00. The hard and soft consistency scores shown in Table 12 are averaged for each metric.

The consistency scores are surprisingly high. All soft consistency scores are above 0.75 means, and the average difference between expert and GPT-4 evaluation in terms of each metric is less than 1 (out of a 5-point scale). The results indicate that GPT-4 might be able to provide a relatively reliable evaluation for machine-generated hypotheses.

### A.9. Evaluation Aspects Description

The evaluation standard for *Validness*, *Novelty*, and *Helpfulness* is correspondingly displayed in Table 13, Table 14, and Table 15.

### A.10. More Details About Past-Feedback Design

In practice, we find that ChatGPT is not capable enough to generate past-feedback with enough good quality for the Inspiration Feedback module. Instead, it tends to provide feedback as “the previous inspiration titles are not very relevant to the hypotheses or the background”. As a result, the ChatGPT Inspiration Title Finder module tends to select inspiration titles that are very related to the background, resulting in a less novel hypotheses generation.

Therefore instead of instantiating with ChatGPT for the Inspiration Feedback module, we experiment with leveraging human

Aspect 1: Helpfulness	
5 points	The hypothesis is novel, valid, clear, and specific enough that it is itself a mature research hypothesis, and human researchers can directly adopt it for publication with no modifications needed.
4 points	The hypothesis is novel enough and can be directly adopted by human researchers for publication after minor modifications.
3 points	The hypothesis should be largely modified or reconstructed by human researchers to adopt it.
2 points	Modifying this hypothesis might not deserve the efforts, but a small part of this hypothesis is inspiring for human researchers to develop a new hypothesis.
1 point	The hypothesis is not helpful and not inspiring at all.

Table 15. Evaluation standard for *Helpfulness*.

heuristics. The heuristics are “if the inspiration titles are less related to the background, then more novel hypotheses are likely to be proposed.”. With this heuristics-based past-feedback, MOOSE does perform better (as shown in the tables in §5 and §6).

This heuristics-based feedback is possible to be obtained by a language model since it has access to the novelty feedback of each hypothesis as well as the inspiration titles the hypothesis leveraged. Here our contribution is to propose a useful framework for the TOMATO task, which is not limited by any LLMs for any module in the framework. In the future, it is possible for more powerful LLMs to find better inspiration feedback than human heuristics.

#### A.11. Qualitative Analysis on Hypotheses Generated From Different Methods

We analyze four hypotheses from the baseline, MOOSE-base, MOOSE-base w/ future-feedback, and MOOSE-base w/ future and past-feedback (MOOSE), where the four methods use the same passage to extract background.

- Hypothesis from the baseline: *Companies that prioritize customer understanding will have higher profitability than companies that do not prioritize customer understanding.* (Evaluated by the expert, Validness: 4; Novelty: 1.5; Helpfulness: 2)
- Hypothesis from MOOSE-base: *The level of empathy displayed by leaders in a startup environment influences employees’ job satisfaction and organizational success through the mediation of employees’ perceived likelihood of negative outcomes and expectation of enjoyment, as well as their propensity toward risky choices.* (Evaluated by the expert, Validness: 3.5; Novelty: 3; Helpfulness: 3)
- Hypothesis from MOOSE-base w/ future-feedback: *Female CMOs in startups, leveraging their higher levels of empathy, are more likely to prioritize customer satisfaction by actively listening to customer feedback, incorporating customer insights into decision-making processes, and providing personalized customer experiences. This employee prioritization of customer satisfaction is positively associated with higher levels of customer engagement, increased customer loyalty, and improved brand recall, as measured by objective metrics such as sales figures, customer retention rates, and brand recognition in the market. The influence of female CMOs on employee behavior is mediated by their ability to foster a caring relationship with customers, as supported by empirical data and statistical analysis.* (Evaluated by the expert, Validness: 3.5; Novelty: 3.5; Helpfulness: 3.5)
- Hypothesis from MOOSE-base w/ future and past-feedback (MOOSE): *Female CMOs’ empathy advantage influences their consideration of negative consequences of CSR initiatives, specifically in terms of employee well-being and job security. This relationship is moderated by individual differences in emotional intelligence. Additionally, the organizational culture and industry context will further influence the relationship between empathy advantage and consideration of negative consequences. The hypothesis will investigate whether female CMOs with higher levels of empathy are more likely to prioritize employee well-being and job security in the implementation of CSR initiatives, and*

*whether this relationship is stronger in industries with a stronger emphasis on employee well-being and job security. It will also explore the mediating role of organizational culture and the moderating role of emotional intelligence in shaping the relationship between empathy advantage and consideration of negative consequences. (Evaluated by the expert, Validness: 4.5; Novelty: 4; Helpfulness: 4)*

Analysis from the expert:

- *H1 falls short of challenging established assumptions or introducing a novel perspective beyond the widely accepted link between customer understanding and profitability.*
- *Both H2 & H3 center around a specific scenario involving female CMOs in startups and delve into their influence on customer satisfaction, employee behavior, and overall business results. From a research standpoint, this more focused approach points to a potential gap in the existing body of knowledge. Moreover, these two hypotheses surpass conventional understanding by considering how the empathy of female CMOs impacts employee behavior and business outcomes. They put forth a fresh viewpoint, suggesting that cultivating a compassionate rapport with customers, fostered by female CMOs, could positively affect customer engagement, loyalty, and brand recognition. These two hypotheses zoom in on a more specific context, introduce an innovative perspective, and probe a potential void in current research. They are anchored in the dynamic world of innovative business settings and propose a more nuanced and all-encompassing connection between variables.*
- *H4 retains its relevance within a modern business landscape by scrutinizing the intersection of empathy, CSR initiatives, and the dynamics of organizations. This syncs seamlessly with the criterion of being rooted in an innovative business environment. Moreover, it shakes up established assumptions by considering the potential adverse outcomes of CSR initiatives and the role empathy plays in shaping decision-making within this context. This hypothesis delves into a more intricate and thorough exploration, examining a broader spectrum of factors and interactions within a specific context. Additionally, it imparts a deeper comprehension of the interplay between empathy, business choices, and organizational results. It grapples with a more complex and distinctive scenario, unearths possible gaps in the existing literature, and introduces a new angle on the role of empathy in the realm of business decisions.*

#### **A.12. Qualitative Analysis on Two MOOSE-Generated Hypotheses With High Expert Evaluation Scores**

In the following two grey boxes are two generated hypotheses from MOOSE with high expert evaluation scores (appended to each hypothesis). The expert's assessment of the two hypotheses is:

*Hypothesis 1: The level of personalization in crowdfunding campaign storytelling, the influence of social media influencers who align with the campaign, the presence of trust indicators, and the emotional appeal of the campaign will positively impact potential donors' likelihood of making a donation. Additionally, the timing of donation requests and the type of social media influencers (e.g., celebrities vs. micro-influencers) will moderate this relationship. The perceived risk associated with the crowdfunding campaign will negatively moderate the relationship between the emotional appeal and donation likelihood. (Validness: 4.5; Novelty: 4.5; Helpfulness: 4.5)*

*Hypothesis 2: Limited financial resources and limited access to networks and markets of women entrepreneurs in the manufacturing sector in developing countries may negatively impact their investment in corporate social responsibility (CSR) initiatives that promote gender equality in host countries. This relationship is further influenced by the intersectionality of gender and race, with women of color facing additional challenges. Additionally, the hypothesis considers the role of institutional factors, such as legal frameworks and policies, and the influence of patriarchal structures on women entrepreneurs' ability to invest in CSR initiatives. (Validness: 3.5; Novelty: 4.0; Helpfulness: 4.0)*

*These two hypotheses both present a comprehensive view of the research narrative. It encompasses multiple hypotheses, including the primary one, as well as the mediation effect, which serves to elucidate the causal connection between the independent and dependent variables. Concurrently, both hypotheses outline the range of the effect — namely, the circumstances in which this effect is applicable, under which scenarios where it might be weakened, and under which situation it could potentially be inverted.*

*In terms of novelty: 1. Limited prior research or a gap in the existing literature. This means that there is a dearth of studies or information available on the subject, making it an unexplored area. 2. Based on a new business setting. It is grounded in an innovative business environment, characterized by novel technologies, contemporary themes, and evolving business requirements. 3. The topic offers a fresh and unique perspective that goes beyond conventional understanding. It might challenge existing assumptions, propose new theories, or present an unconventional approach.*

### A.13. Essential Factors for Good Social Science (and Business) Hypotheses

According to business PhD students, counter-intuitive and novel hypotheses are the mostly favoured (by top business journals). Intuitive and novel hypotheses are also good but not as good as the counter-intuitive ones. Here “novel” refers to “not pointed out by existing literatures”.

Empirically they think of all the hypotheses on top business journals, around 20% are counter-intuitive, leaving the remaining 80% intuitive.

Counter-intuitive hypotheses tend to receive a lower validness evaluation compared to intuitive ones. For this reason, we highlight the counter-intuitive hypothesis in §6.4, even if it receives a lower score in validness than hypotheses in §A.12.

### A.14. An Example of Hypothesis Formulation via MOOSE

Here we show a complete flow of hypothesis discovery, by giving the output of relevant modules. Specifically, the found background, found inspirations, generated suggestions, generated hypothesis, reality feedback, novelty feedback, and clarity feedback are copied in this section.

- Found background:

*Blind boxes, in terms of how they reach the consumers, are an innovation. Digging deeper into the issue, you will learn the very marketing of blind boxes is called probabilistic selling. That is to say, even after making the payment for a blind box, the buyer still has no idea what will be the contents inside. Suspense. A sense of thrill amid an air of expectation ensue. The sheer anticipation of a surprise can be intoxicating even as it can keep the consumer on tenterhooks, in a nice sort of way.*

- Found inspirations:

*Inspiration 0: "Surprise is like crack for your brain. Scientists at Emory and Baylor used MRIs to measure changes in human brain activity in response to a sequence of pleasurable stimuli, using fruit juice and water. The patterns of juice and water squirts were either predictable or completely unpredictable. Contrary to the researchers' expectations, the reward pathways in the brain responded most strongly to the unpredictable sequence of squirts."*

*Inspiration 1: "According to Tmall, China's largest B2C online marketplace, about 200,000 post-95 consumers bought blind boxes from its platform in 2019, with each spending CNY 20,000 per year on average; some fanatics even shelled out CNY 1 million on this "game of chances" to get the coveted limited-edition toys."*

*Inspiration 2: "Product differentiation is fundamentally a marketing strategy to encourage the consumer to choose one brand or product over another in a crowded field of competitors."*

*Inspiration 3: "The pandemic has shown how a lack of solid statistics can be dangerous. But even with the firmest of evidence, we often end up ignoring the facts we don't like"*

*Inspiration 4: "User-generated content (also known as UGC or consumer-generated content) is original, brand-specific content created by customers and published on social media or other channels."*

*Inspiration 5: "Higher Star Ratings Lead to More Orders. More stars really do lead to more sales. The data shows that products with higher star reviews have a much higher volume of purchases."*

- Suggestions from Hypothesis Suggestor Module:

*Suggestion 1: Combining the background information on blind boxes and probabilistic selling with Inspiration 1 (surprise as a powerful reward stimulus) and Inspiration 2 (high spending on blind boxes by young consumers), a novel business research hypothesis could be: The element of surprise and unpredictability in probabilistic selling, such as blind boxes, triggers stronger reward responses in the brain, leading to higher consumer spending and loyalty,*



particularly among younger demographics. This hypothesis suggests that the psychological impact of surprise can be a significant driver of consumer behavior and purchasing decisions, especially for novel product categories.

*Suggestion 2: Integrating the background context with Inspiration 3 (product differentiation strategy) and Inspiration 5 (user-generated content), another innovative hypothesis could be: Encouraging user-generated content around the surprise and reveal aspects of blind box products can serve as a powerful product differentiation strategy, setting brands apart in a competitive market. By leveraging the excitement and social sharing around the unboxing experience, businesses can create a unique brand identity and foster a sense of community among their customers.*

*Suggestion 3: Drawing from the background information, Inspiration 4 (the importance and challenges of evidence-based decision making), and Inspiration 6 (the impact of star ratings on sales), a third hypothesis could be: Despite the proven effectiveness of probabilistic selling in driving consumer engagement and sales, businesses may be reluctant to adopt this strategy due to a lack of traditional, evidence-based metrics and the perceived risks associated with unpredictable outcomes. This hypothesis highlights the potential tension between data-driven decision making and the benefits of innovative, surprise-based selling strategies.*

- Hypothesis from Hypothesis Proposer Module:

*User-generated content, such as unboxing videos and reviews, can significantly influence the purchasing decisions of potential blind box buyers, particularly among younger demographics. The excitement and social sharing around the unboxing experience, as depicted in user-generated content, can amplify the psychological reward of surprise and drive higher sales, making it an effective tool for businesses to differentiate their products and increase consumer engagement in the blind box market.*

- Reality Feedback:

*The given hypothesis about the influence of user-generated content on blind box purchasing decisions among younger demographics seems to reflect reality to a certain extent. However, there are a few points to consider that might suggest the hypothesis does not entirely capture the reality:*

*1. Generalization: The hypothesis assumes that user-generated content significantly influences all younger demographics. In reality, the impact may vary depending on factors such as personal preferences, cultural background, and socio-economic status.*

*2. Other influential factors: While user-generated content can be an important factor in purchasing decisions, it may not be the sole or most significant driver. Other aspects like product quality, price, brand reputation, and peer influence might play equally important or even more crucial roles.*

*3. Limited scope: The hypothesis focuses on the excitement and social sharing aspect of unboxing experiences. However, it does not consider the potential negative impact of user-generated content, such as disappointing unboxing experiences or negative reviews that could discourage purchases.*

*4. Assumption of psychological reward: The hypothesis assumes that the surprise element in blind box purchases always leads to a positive psychological reward. In reality, some consumers might experience disappointment or frustration if they do not receive the desired item, which could deter future purchases.*

*5. Lack of empirical evidence: The hypothesis presents a plausible explanation but does not provide concrete empirical evidence*

- Novelty Feedback:

*Based on the provided inspirations and the possibly related business literature paragraph, the research hypothesis appears to be novel and not a direct semantic copy of any of the given information. The hypothesis combines insights from various sources, such as the psychological impact of surprise, the popularity of blind boxes among younger demographics, the role of product differentiation, and the influence of user-generated content on purchasing decisions. While the hypothesis draws from these inspirations, it synthesizes them into a unique and specific context related to the blind box market and the impact of user-generated content on consumer engagement and sales. The provided literature paragraph discusses language modeling and large language models, which is not directly related to the hypothesis. To further enhance the novelty of the hypothesis, you could consider the following suggestions:*

*1. Explore the potential long-term effects of user-generated content on brand loyalty and customer retention in the blind box market.*

2. Investigate the role of social media platforms and influencers in amplifying the impact of user-generated content on blind box sales.
3. Examine the potential differences in the effectiveness of user-generated content across various product categories within the blind box market (e.g., collectible toys, beauty products, or snacks).
4. Analyze the impact of user-generated content on the perceived value and collectibility of blind box items among consumers.

- **Clarity Feedback:**

*The given hypothesis is relatively clear and specific, as it identifies the relationship between user-generated content and purchasing decisions of potential blind box buyers. It also specifies the target demographic (younger buyers), the psychological factors at play (excitement, surprise, and social sharing), and the potential benefits for businesses (product differentiation and increased consumer engagement). However, there is still room for improvement in terms of specificity. To make the hypothesis more specific, consider the following:*

1. *Quantify the extent of the influence of user-generated content on purchasing decisions (e.g., "can increase the likelihood of purchase by X").*
2. *Specify the age range of the "younger demographics" (e.g., "particularly among buyers aged 18-34").*
3. *Clarify which types of user-generated content have the most significant impact (e.g., "especially unboxing videos that showcase the excitement and surprise of the experience").*
4. *Indicate the platforms or channels where user-generated content is most effective (e.g., "on social media platforms such as Instagram and TikTok").*
5. *Provide a more detailed explanation of how user-generated content amplifies the psychological reward (e.g., "by creating a sense of anticipation and vicarious experience for potential buyers").*
6. *Quantify the potential*

#### **A.15. Future Directions**

This work discovered the possibility of LLMs to propose novel research hypotheses. But it mainly focuses on the social science and business disciplines. It would be very interesting to investigate how LLMs can induce novel hypotheses for other disciplines (especially nature science domains).

In addition, the MOOSE framework could be further improved to induce more valid and novel hypotheses for social science and business domains.

From the aspect of human-AI interaction, it would be also interesting to see how MOOSE can act as an AI Copilot to assist scientists in hypothesis discovery.

#### **A.16. Full Algorithm of the MOOSE Framework**

Algorithm 1 shows the full algorithm of the proposed framework.

---

**Algorithm 1** Algorithm for MOOSE

---

**Input:** Raw web corpus  $C$ , related surveys  $S$ **Parameter:** Total iterations for *past-feedback*  $M$ , total iterations for *present-feedback*  $N$ **Output:** A list of hypotheses  $H$ 

```
1: for  $c$  in  $C$  do
2:    $b, b\_reason = \text{Background\_Finder}(c)$ 
3:   if  $b == \text{None}$  then
4:     continue
5:   end if
6:   for iteration  $k \in 0 \dots M$  do
7:     if  $k \neq 0$  then
8:        $past\_f = \text{Inspiration\_Feedback}(t, h, present\_f)$ 
9:     else
10:       $past\_f = \text{None}$ 
11:    end if
12:     $t, t\_reason = \text{Inspiration\_Title\_Finder}(C, b, b\_reason, past\_f)$ 
13:     $p = \text{find\_passage\_by\_title}(t, C)$ 
14:     $i = \text{Inspiration\_Finder}(b, b\_reason, p, t\_reason)$ 
15:     $s = \text{Hypothesis\_Suggestor}(b, i)$ 
16:     $h = \text{Hypothesis\_Proposer}(b, i, s)$ 
17:    for iteration  $t \in 0 \dots N$  do
18:       $cf, rf, nf = \text{Clarity\_Checker}(h), \text{Reality\_Checker}(h), \text{Novelty\_Checker}(h, S)$ 
19:       $present\_f = [cf, rf, nf]$ 
20:       $h = \text{Hypothesis\_Proposer}(b, i, s, h, present\_f)$ 
21:    end for
22:     $H.append(h)$ 
23:  end for
24: end for
25: return  $H$ 
```

---