Large Language Models for Automated Open-domain Scientific Hypotheses Discovery

Anonymous ACL submission

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 hand-800 picked 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 011 012 tackle these problems by proposing the first NLP dataset for social science academic hypotheses discovery, consisting of 50 recent top 014 social science publications; and a raw web cor-015 pus that contains enough information to make 017 it possible to develop all the research hypotheses in the 50 papers. The final goal is to create systems that automatically generate valid, novel, and helpful scientific hypotheses, given only a pile of raw web corpus. Different from the 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, as well as three different feedback mechanisms that em-027 pirically show performance gain over the base framework. Finally, our framework exhibits superior performance in terms of both GPT-4 based evaluation and expert-based evaluation.

1 Introduction

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



Figure 1: Overview of the new task setting of hypothetical induction and the role of the MOOSE framework.

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.

044

045

046

047

048

051

053

054

056

057

059

060

061

062

063

064

065

066

067

068

069

070

071

072

The latest research on hypothetical induction (Yang et al., 2022b) has two main limitations. 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)¹. This hypothesis formation process is seen as the first step for scientific discovery (Wang et al., 2023a). We call this task as "au**TO**mated open-do**MA**in hypo**T**hetical inducti**O**n (TOMATO)". It is "automated" since a method for this task should automatically propose

¹Dataset and Code available at https://anonymous. 4open.science/r/TOMATO/.

hypotheses with few human efforts; It is opendomain since it is not restricted by any manually 074 collected data. For the TOMATO task, we con-075 structed 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 079 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 multimodule 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.

097

099

101

102

103

104

105

106

108

109

110

111

112

113

114

115

116

117

118

119

120

121

122

123

For both GPT-4 (OpenAI, 2023) evaluation and social science expert evaluation, our experiment indicates that our framework performs better than an LLM (Ouyang et al., 2022) based baseline, and124each of the three feedback mechanisms can progressively improve the base framework. During human125analysis, many hypotheses generated by our framework are recognized by social science researchers127to be valid, novel, and helpful in the same time.128

130

131

132

133

134

135

136

137

138

139

140

141

142

143

144

145

146

147

148

149

150

151

152

153

154

155

156

157

158

159

162

163

164

165

166

168

169

170

171

172

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 language-described 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 ground truth discovery is not from recent publications. Wang et al. (2023b) propose an automatic method to collect NLP publications to construct a dataset, and a method to propose hypotheses in the NLP domain. However, (1) their task needs human-provided input, while our task/method support both human-provided input and automatically searched input; (2) their dataset is not manually collected, and their background text and seed terms are collected in the same paper which proposes the ground truth hypothesis, which might cause data contamination problem; (3) their dataset is composed of ACL anthology papers before 2021, so the papers in the dataset are likely to appear in the training corpus of ChatGPT as well as LLaMA-based models (Touvron et al., 2023); (4) their method does not leverage feedback mechanism and is not specifically designed to propose novel hypotheses. 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 futurefeedback), 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., 2022a; 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 **Hypothesis 2.** Customers whose preceding customers use FR payment technology are more likely to use FR payment technology than those whose 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.

over-generate-then-filter mechanism (Yang et al., 2022b). However, they only utilize LLMs to "filter" but not to provide feedback.

3 Dataset Collection

173

174

175

176

177

178

180

181

184

188

190

192

195

196

197

198

201

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 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

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

Table 1: Statistics of subject distribution of the dataset.

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). 204

205

206

207

208

209

210

211

212

213

214

215

216

217

218

219

220

221

222

223

226

227

229

230

231

234

235

236

237

238

239

240

241

242

243

244

246

247

248

249

250

251

Instead of directly copying the background and inspirations 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 multimodule framework and three feedback mechanisms (past-feedback, present-feedback, and futurefeedback). We call the full framework as **M**ultim**O**dule framew**O**rk with pa**S**t present future fe**E**dback (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 soci-



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.

etal 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 viable 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 $M_0, M_1, ..., M_n$ 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 .

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., 2022b) 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. 287

289

290

291

292

293

294

295

296

297

298

299

300

301

302

303

304

305

306

307

308

309

310

311

312

313

314

315

316

317

318

319

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. (2022b) 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

324

325

326

330

331

334

336

341

342

343

346

347

351

358

367

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 an additional module utilizes o_i , o_j , and f_j 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. 371

372

373

374

375

376

377

379

381

384

385

386

388

390

391

392

393

394

395

396

397

398

399

400

401

402

403

404

405

406

407

408

409

410

411

412

413

414

415

416

417

418

419

420

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., 2022b). 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.

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 and 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 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

	Validness	Novelty	Helpfulness
Baseline	3.954	2.483	3.489
MOOSE-base	3.907	3.081	3.859
w/ future-feedback	3.955	3.226	3.953
w/ future- and past-feedback	3.916	3.390 †	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*. Results with [†] mean the difference compared to the baseline is statistically significant (p < 0.01) using Bootstrap method (Berg-Kirkpatrick et al., 2012).

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

Table 3: Effect of *present-feedback* (evaluated by *GPT-*4). Results with [†] mean the difference compared to MOOSE w/o *present-feedback* is significant (p < 0.01).

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 machinegenerated social science hypotheses.

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.

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.

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

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

Table 4: Effect of MOOSE-base, *future-feedback* and *past-feedback* (evaluated by *experts*). MOOSE results are selected from the 5th iteration of *present-feedback*.

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*4*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.

450

451

452

453

454

455

456

457

458

459

460

461

462

463

464

465

466

467

468

469

470

471

472

473

474

475

476

477

478

479

480

481

482

483

484

485

486

487

488

489

490

491

Table 2 shows GPT-4's evaluation targeting at comparing MOOSE-base and the baseline and shows the effect of future-feedback and pastfeedback. In this table, MOOSE-related results are averaged over iterations of present-feedback to not be influenced by present-feedback. MOOSEbase 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. MOOSEbase with future-feedback comprehensively outperforms MOOSE-base in terms of all three metrics. MOOSE-base with both future and pastfeedback 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 presentfeedback 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 MOOSEbase with both future and past-feedback to focus on present-feedback. It shows that as more iterations of present-feedback are conducted, validness and

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

Table 5: Effect of present-feedback (eval. by experts).

	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	3.364	3.670
GPT-3.5 picked background and inspirations	3.812	2.818	3.733
Groundtruth background and inspirations	3.876	3.000	3.806
Groundtruth hypotheses	3.700	3.380	3.880

Table 6: Analysis of retrieval's effect on generated hypotheses (evaluated by *GPT-4*). No methods here utilize any feedback mechanisms. Every method here uses the same ChatGPT-based hypothesis proposer module.

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 pastfeedback. MOOSE-related results are selected from the 5^{th} 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 presentfeedback. 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.

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

	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	3.431	3.886
w/ randomized corpus	3.941	3.227	3.955

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

524

525

526

527

528

529

530

531

532

533

534

535

536

537

538

539

540

541

542

543

544

545

546

547

548

549

550

551

552

553

554

555

556

557

558

559

560

561

562

563

564

565

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, Chat-GPT 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 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. In-

518

519

520

521

522

523

581

584

585

586

587

588

589

591

595

599

600

Here is the assessment from one of the experts:

stead, FF2 (the hypothesis suggestor module) has 566 the global view and only provides soft initial sug-567 gestions 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 571 and other inspirations to propose hypotheses. To be fair, MOOSE-base, which is not equipped with the 573 hypothesis suggestor module, has the same prompt to combine the inspirations together (just without 575 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 presentfeedback, 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 Qualitative Analysis

The following box shows one generated counterintuitive 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) The main reason I give a high mark for both three dimensions of this hypothesis is because:

610

611

612

613

614

615

616

617

618

619

620

621

622

623

624

625

626

627

628

629

630

631

632

633

634

635

636

637

638

639

640

641

642

643

644

645

646

647

648

649

650

651

652

653

654

655

656

657

658

659

(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 individualism that have been widely studied in cross-cultural psychology.

(2) For novelty, this hypothesis is counterintuitive 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 wellbeing 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.

More analyses are in A.11 (hypotheses comparison), A.12 (high expert evaluation examples), and A.13 (factors for good hypotheses).

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 MOOSEbase 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 high-quality research hypotheses.

Limitations

660

662

670

671

674

675

676

678

702

704

706

707

708

710

From the first look, it might seem that the proposed dataset consists of only 50 recent papers. However, they are all manually collected by experts (PhD students), and are annotated with lots of details (e.g., identifying background and inspirations, finding relevant raw web passages for background and inspirations, reasoning process, complexity level). In addition, each paper has been published in a top social science journal, representing the pinnacle of human intelligence. This means it would be incredibly exciting if LLMs could propose a hypothesis from even a single one of these recent papers.

It might also seem that it is not clear whether the design of the framework can apply to other disciplines. However, to the best of our knowledge, this is the first paper using LLMs that can propose novel scientific hypotheses that are new to humanity. We choose social science as the breakthrough point since the main data format of social science is language. Table 1 shows that the dataset covers 7 different disciplines (e.g., Psychology, Management, Marketing). It would be nearly impossible for the first few works to develop a general method to propose novel hypotheses for all disciplines.

This paper concentrates on an automated task setting in which a system is designed to formulate scientific hypotheses independently, without requiring human intervention. In some scenarios, scientists may prefer to use their own background and inspirations as input for controllable hypotheses generation. It might seem that the automated setting and the controllable setting are in conflict. However, we contend that the automated setting make a step further than the controllable setting, since a system developed for an automated setting would inherently support controllable generation by simply substituting the automatically searched inputs (e.g., background and inspirations) with those that are manually crafted.

Societal Impact: Expert evaluation shows that MOOSE, an LLM-based system, might already be able to serve as a copilot for researchers across various social science disciplines. Particularly, as depicted in Figure 1, it can assist researchers in the hypothesis formation process, which is the first step for scientific discovery (Wang et al., 2023a). This capability signifies a step towards enhancing the efficiency of scientific exploration by accelerating the formation and development of innovative and credible research hypotheses, thereby boosting researchers' productivity. To maximize its impact and ensure equitable access, it is imperative to advocate for the open-sourcing of such systems, thereby democratizing access for the global scientific community. 711

712

713

714

715

716

717

718

719

720

721

722

723

724

725

726

727

728

729

730

731

732

733

734

735

736

738

739

740

741

742

743

744

745

746

747

748

749

750

751

752

753

754

755 756

757

758

759

760

761

762

763

Ethics Statement

This article follows the ACL Code of Ethics. To our knowledge, there are no foreseeable potential risks to use the dataset and methods in this paper.

References

- Satanjeev Banerjee and Alon Lavie. 2005. 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*, pages 65–72, Ann Arbor, Michigan. Association for Computational Linguistics.
- Taylor Berg-Kirkpatrick, David Burkett, and Dan Klein. 2012. An empirical investigation of statistical significance in nlp. In *Proceedings of the 2012 Joint Conference on Empirical Methods in Natural Language Processing and Computational Natural Language Learning*, pages 995–1005.
- Daniil A. Boiko, Robert MacKnight, and Gabe Gomes. 2023. Emergent autonomous scientific research capabilities of large language models. *CoRR*, abs/2304.05332.
- Antoine Bosselut, Hannah Rashkin, Maarten Sap, Chaitanya Malaviya, Asli Celikyilmaz, and Yejin Choi. 2019. COMET: Commonsense transformers for automatic knowledge graph construction. In *Proceedings* of the 57th Annual Meeting of the Association for Computational Linguistics, pages 4762–4779, Florence, Italy. Association for Computational Linguistics.
- Andres M Bran, Sam Cox, Andrew D White, and Philippe Schwaller. 2023. Chemcrow: Augmenting large-language models with chemistry tools. *arXiv preprint arXiv:2304.05376*.
- Rajarshi Das, Manzil Zaheer, Dung Thai, Ameya Godbole, Ethan Perez, Jay Yoon Lee, Lizhen Tan, Lazaros Polymenakos, and Andrew McCallum. 2021. Casebased reasoning for natural language queries over knowledge bases. In Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing, pages 9594–9611, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Jia Gao, Ying Rong, Xin Tian, and Yuliang Yao. 2023. Improving convenience or saving face? an empirical analysis of the use of facial recognition payment technology in retail. *Information Systems Research*.

- 765 766 767
- 76 76

- 7
- 775

777 778

779 780

781 782

7

784 785

7

787 788

- 8
- 8

809 810

811 812

813 814

815

81

817 818

- Vinod Goel, Gorka Navarrete, Ira A Noveck, and Jérôme Prado. 2017. The reasoning brain: The interplay between cognitive neuroscience and theories of reasoning.
- Chin-Yew Lin. 2004. Rouge: A package for automatic evaluation of summaries. In *Text summarization branches out*, pages 74–81.
- Chia-Wei Liu, Ryan Lowe, Iulian Serban, Mike Noseworthy, Laurent Charlin, and Joelle Pineau. 2016.
 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*, pages 2122–2132, Austin, Texas. Association for Computational Linguistics.
 - Aman Madaan, Niket Tandon, Prakhar Gupta, Skyler Hallinan, Luyu Gao, Sarah Wiegreffe, Uri Alon, Nouha Dziri, Shrimai Prabhumoye, Yiming Yang, Sean Welleck, Bodhisattwa Prasad Majumder, Shashank Gupta, Amir Yazdanbakhsh, and Peter Clark. 2023. Self-refine: Iterative refinement with self-feedback. *CoRR*, abs/2303.17651.

John D Norton. 2003. A little survey of induction.

- OpenAI. 2023. GPT-4 technical report. CoRR, abs/2303.08774.
- Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll L. Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser Kelton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul F. Christiano, Jan Leike, and Ryan Lowe. 2022. Training language models to follow instructions with human feedback. In *NeurIPS*.
- Feng Pan, Rutu Mulkar-Mehta, and Jerry R. Hobbs. 2011. Annotating and learning event durations in text. *Comput. Linguistics*, 37(4):727–752.
- Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the* 40th annual meeting of the Association for Computational Linguistics, pages 311–318.
- Baolin Peng, Michel Galley, Pengcheng He, Hao Cheng, Yujia Xie, Yu Hu, Qiuyuan Huang, Lars Liden, Zhou Yu, Weizhu Chen, and Jianfeng Gao. 2023. Check your facts and try again: Improving large language models with external knowledge and automated feedback. *CoRR*, abs/2302.12813.
- Ofir Press, Muru Zhang, Sewon Min, Ludwig Schmidt, Noah A. Smith, and Mike Lewis. 2022. Measuring and narrowing the compositionality gap in language models. *CoRR*, abs/2210.03350.
- Biqing Qi, Kaiyan Zhang, Haoxiang Li, Kai Tian, Sihang Zeng, Zhang-Ren Chen, and Bowen Zhou. 2023.
 Large language models are zero shot hypothesis proposers. *CoRR*, abs/2311.05965.

Stephen Robertson, Hugo Zaragoza, et al. 2009. The probabilistic relevance framework: Bm25 and beyond. *Foundations and Trends*® *in Information Retrieval*, 3(4):333–389. 819

820

821

822

823

824

825

826

827

828

829

830

831

832

833

834

835

836

837

838

839

840

841

842

843

844

845

846

847

848

849

850

851

852

853

854

855

856

857

858

859

860

861

862

863

864

865

866

867

868

869

870

871

873

- Noah Shinn, Beck Labash, and Ashwin Gopinath. 2023. Reflexion: an autonomous agent with dynamic memory and self-reflection. *CoRR*, abs/2303.11366.
- Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurélien Rodriguez, Armand Joulin, Edouard Grave, and Guillaume Lample. 2023. Llama: Open and efficient foundation language models. *CoRR*, abs/2302.13971.
- Hanchen Wang, Tianfan Fu, Yuanqi Du, Wenhao Gao, Kexin Huang, Ziming Liu, Payal Chandak, Shengchao Liu, Peter Van Katwyk, Andreea Deac, et al. 2023a. Scientific discovery in the age of artificial intelligence. *Nature*, 620(7972):47–60.
- Qingyun Wang, Doug Downey, Heng Ji, and Tom Hope. 2023b. Learning to generate novel scientific directions with contextualized literature-based discovery. *CoRR*, abs/2305.14259.
- Kevin Yang, Yuandong Tian, Nanyun Peng, and Dan Klein. 2022a. Re3: Generating longer stories with recursive reprompting and revision. In *Proceedings* of the 2022 Conference on Empirical Methods in Natural Language Processing, EMNLP 2022, Abu Dhabi, United Arab Emirates, December 7-11, 2022, pages 4393–4479. Association for Computational Linguistics.
- Zonglin Yang, Li Dong, Xinya Du, Hao Cheng, Erik Cambria, Xiaodong Liu, Jianfeng Gao, and Furu Wei. 2022b. Language models as inductive reasoners. *CoRR*, abs/2212.10923.
- Zonglin Yang, Xinya Du, Erik Cambria, and Claire Cardie. 2023a. 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, pages 3509–3522, Dubrovnik, Croatia. Association for Computational Linguistics.
- Zonglin Yang, Xinya Du, Rui Mao, Jinjie Ni, and Erik Cambria. 2023b. Logical reasoning over natural language as knowledge representation: A survey. *CoRR*, abs/2303.12023.
- Zonglin Yang, Xinya Du, Alexander Rush, and Claire Cardie. 2020. Improving event duration prediction via time-aware pre-training. In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 3370–3378, Online. Association for Computational Linguistics.
- Ruiqi Zhong, Peter Zhang, Steve Li, Jinwoo Ahn, Dan Klein, and Jacob Steinhardt. 2023. Goal driven discovery of distributional differences via language descriptions. *CoRR*, abs/2302.14233.

966

A Appendix

875

884

894

895

900

901

902

903

904

905

906

907

908

909

910

911

912

913

914

915

916

A.1 Hyper-parameters, Anonymous Github Link, and Full Prompts

All experiments are conducted with gpt-3.5-turbo, with 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.

The dataset and code of this submission are already public on GitHub. An anonymous version can be found at: https: //anonymous.4open.science/r/TOMATO/.

Particularly, the dataset can be found at https://anonymous.4open.science/r/

TOMATO/Data/business_research.xlsx.

The full prompts for MOOSE framework is shown in prompts_for_tomato_modules() function in utils.py.

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.

Qi et al. (2023) work on zero-shot hypothesis proposing, which is a concurrent work to our paper. They don't focus on social science and business disciplines, and mostly adopt a single LLM as method (prompting, finetuning).

A.3 Dataset Complexity Distribution

Table 8 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

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

Table 8: Statistics of the complexity of the dataset.

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 and 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 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).

967

968

969

970

972

973 974

975

978

979

982

987

991

993

996

997

1002

1003

1004

1005

1006

1007

1008

1010

1011

1012

1013

1014

1015

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) MOOSEbase + future-feedback; (4) MOOSE-base + futurefeedback + past-feedback. For (2) and (4), we collect three hypotheses, which are (a) without presentfeedback; (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

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

Table 9: 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 10: Hard and soft consistency scores between expert evaluation and GPT-4 evaluation in terms of Validness, Novelty, and Helpfulness metrics.

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. 1016

1017

1018

1019

1020

1021

1022

1023

1024

1025

1026

1027

1029

1030

1031

1032

1033

1034

1035

1036

1037

1038

1039

1040

1041

1042

1043

1044

1045

1046

1047

1048

1049

The consistency scores between experts are shown in Table 9. 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.

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 10 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

1051

1088 1089

1090 1091

1093 1094

1095

1096

1097 1098 1099

1100

evaluation for machine-generated hypotheses.

A.9 Evaluation Aspects Description

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.

Aspect 2: Novelty.

5 points: the hypothesis is completely novel and has not been proposed by any existing literature;

4 points: the main argument or several subarguments 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.

Aspect 3: 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.

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.

1101

1102

1103

1104

1105

1106

1107

1108

1109

1110

1111

1112

1113

1114

1115

1116

1117

1118

1119

1120

1121

1122

1123

1124

1125

1126

1127

1128

1129

1130

1131

1132

1133

1134

1135

1136

1137

1138

1139

1140

1141

1142

1143

1144

1145

Therefore instead of instantiating with ChatGPT for the Inspiration Feedback module, we experiment with leveraging human 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 $\S5$ and $\S6$).

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.

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

We analyze four hypotheses from the baseline, MOOSE-base, MOOSE-base w/ futurefeedback, and MOOSE-base w/ future and pastfeedback (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-1146 feedback: Female CMOs in startups, lever-1147 aging their higher levels of empathy, are more 1148 likely to prioritize customer satisfaction by 1149

actively listening to customer feedback, in-1150 corporating customer insights into decision-1151 making processes, and providing personalized 1152 customer experiences. This employee priori-1153 tization of customer satisfaction is positively 1154 associated with higher levels of customer en-1155 gagement, increased customer loyalty, and 1156 improved brand recall, as measured by objec-1157 tive metrics such as sales figures, customer 1158 retention rates, and brand recognition in the 1159 market. The influence of female CMOs on em-1160 ployee behavior is mediated by their ability 1161 to foster a caring relationship with customers, 1162 as supported by empirical data and statistical 1163 analysis. (Evaluated by the expert, Validness: 1164 3.5; Novelty: 3.5; Helpfulness: 3.5) 1165

• Hypothesis from MOOSE-base w/ future and 1166 past-feedback (MOOSE): Female CMOs' em-1167 pathy advantage influences their considera-1168 tion of negative consequences of CSR initia-1169 tives, specifically in terms of employee well-1170 being and job security. This relationship is 1171 moderated by individual differences in emo-1172 tional intelligence. Additionally, the organi-1173 zational culture and industry context will fur-1174 ther influence the relationship between em-1175 pathy advantage and consideration of nega-1176 tive consequences. The hypothesis will investi-1177 gate whether female CMOs with higher levels 1178 of empathy are more likely to prioritize em-1179 ployee well-being and job security in the im-1180 plementation of CSR initiatives, and whether 1181 this relationship is stronger in industries with 1182 a stronger emphasis on employee well-being 1183 and job security. It will also explore the me-1184 diating role of organizational culture and the 1185 moderating role of emotional intelligence in 1186 shaping the relationship between empathy ad-1187 vantage and consideration of negative conse-1188 quences. (Evaluated by the expert, Validness: 1189 4.5; Novelty: 4; Helpfulness: 4) 1190

Analysis from the expert:

1191

1192

1193

1194

1195

1196

1197

1198

1199

- 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 1200 focused approach points to a potential gap 1201 in the existing body of knowledge. Moreover, 1202 these two hypotheses surpass conventional un-1203 derstanding by considering how the empathy 1204 of female CMOs impacts employee behavior 1205 and business outcomes. They put forth a fresh 1206 viewpoint, suggesting that cultivating a com-1207 passionate rapport with customers, fostered 1208 by female CMOs, could positively affect cus-1209 tomer engagement, loyalty, and brand recogni-1210 tion. These two hypotheses zoom in on a more 1211 specific context, introduce an innovative per-1212 spective, and probe a potential void in current 1213 research. They are anchored in the dynamic 1214 world of innovative business settings and pro-1215 pose a more nuanced and all-encompassing 1216 connection between variables. 1217

• H4 retains its relevance within a modern busi-1218 ness landscape by scrutinizing the intersection 1219 of empathy, CSR initiatives, and the dynam-1220 ics of organizations. This syncs seamlessly 1221 with the criterion of being rooted in an in-1222 novative business environment. Moreover, it 1223 shakes up established assumptions by consid-1224 ering the potential adverse outcomes of CSR 1225 initiatives and the role empathy plays in shap-1226 ing decision-making within this context. This 1227 hypothesis delves into a more intricate and 1228 thorough exploration, examining a broader 1229 spectrum of factors and interactions within 1230 a specific context. Additionally, it imparts 1231 a deeper comprehension of the interplay be-1232 tween empathy, business choices, and orga-1233 nizational results. It grapples with a more 1234 complex and distinctive scenario, unearths 1235 possible gaps in the existing literature, and 1236 introduces a new angle on the role of empathy 1237 in the realm of business decisions. 1238

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

1239

1240

1241

1242

1243

1244

1245

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:

These two hypotheses both present a comprehen-
sive view of the research narrative. It encompasses1246multiple hypotheses, including the primary one, as1247well as the mediation effect, which serves to elu-1249

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. microinfluencers) 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; Helpfulness: 4)

cidate 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.

1250

1251

1252

1253

1254

1255

1256

1257

1258

1259

1260

1261

1263

1264

1265

1266

1267

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.

1268

1269

1270

1271

1272

1273

1274

1275

1276

1277

1278

1279

1281

1282

1283

1284

1285

1286

1287

1288

1289

1290

1291

1292

1293

1294

1295

1296

1297

1298

1299

1300

1301

1303

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

According to social science PhD students, counterintuitive and novel hypotheses are the mostly favoured (by top social science and 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 social science journals, around 20% are counterintuitive, 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 §6.3, even if it receives a lower score in validness than hypotheses in §A.12.

A.14 Prompts for Each Modules in MOOSE

This section illustrates the prompt MOOSE adopts.

Background Finder: In the provided passage, likely from a business-related report, try to collect the best paragraph (or sentence) in the reports that could serve as suitable academic background for social science research. The chosen academic background in business should encompass research topics that can be further developed into hypotheses for social science research. The passage is: passage. Please give a response to the initial question of exactly extracting the best business academic background paragraph (or sentence) from the given passage, and also provide an evaluation of the selected background in terms of what are possible social science research directions given the background (response format: 'Background: Evaluation: ...').

Inspiration Title Finder: Given an academic 1304 background in social science research and titles 1305 of business-related reports, which titles (and their 1306 corresponding business reports) could contain re-1307 search inspirations which combined with the back-1308 ground could lead to non-trivial hypotheses in 1309 social science research? The academic back-1310 ground is *background*. The title collections are 1311 *title_collections*. Please give a response to the 1312 initial question of extracting three titles that most 1313 probably contain suitable research inspirations 1314 given the social science research background, and 1315 also evaluate the selected titles in terms of how it 1316

could potentially help social science research hypothesis developing (response format: 'Title: Evaluation: Title: Evaluation: ...').

Inspiration Finder: Given an academic back-1320 ground in social science research and a business-1321 related report, try to collect the best one sentence 1322 1323 or one paragraph in the report that possibly contain an inspiration, which could be used together 1324 with the given background to further develop a 1325 hypothesis in social science research (usually a hypothesis is more novel if its inspiration is less 1327 directly related to the given background). The aca-1328 demic background is *background*. The business 1329 report is: report. Previous feedback on how this 1330 1331 passage could possibly contribute to a hypothesis by only seeing the title of this inspiration passage: 1332 *feedback*. Please give a response to the initial 1333 question of exactly extracting the best one sentence 1334 or one paragraph from the business-related report 1335 (but not from background or evaluation of titles) 1336 as a possible inspiration, and also evaluate the ex-1337 tracted inspiration in terms of its own quality, how 1338 it can potentially help social science research hy-1339 pothesis developing, and how is it related to given 1340 background (response format: 'Inspiration: Evalu-1341 ation: '). 1342

Inspiration Feedback: Given an academic back-1343 ground in social science research, previously selected titles of business-related reports, previously 1345 generated social science research hypothesis using 1346 the academic background and some inspirations 1347 from the selected reports (according to selected titles for reports), and evaluation of previously 1349 generated hypothesis, try to understand potential 1350 problems of previously generated social science 1351 research hypothesis that might be caused by im-1352 proper selection of business reports, and identify 1353 potential problems of report selection. The aca-1354 demic background is *background*. The previously 1355 selected titles are: *titles*. The previously generated hypotheses and their evaluation are: hypotheses 1357 and *feedback*. Please give a response to the initial 1358 question of identifying and elaborating problems 1359 of the previously selected report titles that might cause negative effect on generating the given spe-1361 cific hypothesis. 1362

Hypotheses Suggestor:Given an academic back-ground in social science research and some possibleinspirations which combined with the backgroundcould lead to meaningful social science research

hypothesis, please try to give some suggestions 1367 on how these inspirations could be combined to 1368 be potentially helpful to propose novel social sci-1369 ence research hypotheses. Multiple inspirations 1370 are encouraged to be used together to generate 1371 new hypotheses. Inspirations which seem to be 1372 less connected to the background could probably 1373 contribute more to a novel hypothesis. A good 1374 business hypothesis should be novel and not intu-1375 itive, should has never been formally proposed in 1376 the social science research fields ever before. The 1377 background is: background. The possible inspi-1378 rations are: *inpirations*. Please give a response 1379 to the initial question of generating suggestions 1380 on how the background and inspirations could be 1381 combined to generate novel social science research 1382 hypotheses. Each suggestion should leverage more 1383 than two inspirations (response format: 'Sugges-1384 tion 1: Suggestion 2: ...') 1385

Hypotheses Proposer: Given an academic back-1386 ground in social science research, some possible 1387 inspirations which combined with the background 1388 could lead to meaningful social science research 1389 hypothesis, and some initial suggestions on how to 1390 leverage these inspirations to build hypotheses, try to give unique hypotheses based on the background, 1392 inspirations, and the initial suggestions. Multiple 1393 inspirations and suggestions are encouraged to be 1394 used together to generate new hypotheses. Inspira-1395 tions which seem to be less connected to the back-1396 ground could probably contribute more to a novel 1397 hypothesis. A good business hypothesis should (1) 1398 contain an independent variable and a dependent 1399 variable, and describe how the independent variable 1400 can influence the dependent variable, and (2) be 1401 novel and not intuitive, should has never been for-1402 mally proposed in the social science research fields 1403 ever before. The background is: background. The 1404 possible inspirations are: inspirations. The sug-1405 gestions are: *suggestions*. Please give a response 1406 to the initial question of generating unique mean-1407 ingful social science research hypotheses given the 1408 background, inspirations, and suggestions. Each 1409 hypothesis should leverage more than two sugges-1410 tions or inspirations. For each hypothesis, please 1411 give the reasoning processing first, and then give 1412 the hypothesis. (response format: 'Reasoning pro-1413 cess: Hypothesis: Reasoning process: Hypothesis: 1414 ...'). 1415

Reality Checker: Given a research hypothesis in 1416

1417social science research, try to give some feedback1418on whether the hypothesis by any chance does not1419reflects the reality. Please directly answer this ques-1420tion. The hypothesis is: hypothesis. Please give a1421response to the initial question of providing feed-1422back on whether the research hypothesis reflects1423the reality.

Novelty Checker: Given a research hypothesis in 1424 social science research, some inspirations used for 1425 developing the hypothesis, and a possibly related 1426 paragraph from a relevant social science research 1427 survey, try to give some feedback on whether the 1428 hypothesis is by any chance not novel (the reason 1429 is that the hypothesis is used for social science 1430 research, where novel and not ever proposed hy-1431 potheses are preferred). To be novel, the hypoth-1432 esis should at least not be semantically a direct 1433 copy of any inspiration or any arguments in ex-1434 isting business literature (including literature that 1435 are not provided as input), but could be a conclu-1436 sion from multiple reasoning steps using the inspi-1437 1438 rations, and probably then with (slightly / some) deviations from the conclusion. The hypothesis is: 1439 hypothesis. The inspirations used for developing 1440 the hypothesis are: inspirations. One of the most 1441 similar existing business literature paragraph is: 1442 *paragraph*. Please give a responses to the initial 1443 1444 question of providing detailed feedback on whether the research hypothesis is by any chance not novel 1445 (not a semantically direct copy of any inspiration 1446 or any argument in existing business literature). 1447

Clarity Checker: Given a research hypothesis in social science research, try to give some feedback on whether the hypothesis is clear and specific enough. By specific, it means a hypothesis should not only indicate two elements are related, but also how they are related, to what extent they are related, why they are related, and which specific sub-elements of the two elements are related. The hypothesis is: *hypothesis*. Please give a response to the initial question on whether the hypothesis is clear and specific enough.

A.15 Full Algorithm for the Proposed Multi-Module Framework

Algorithm 1 shows the full algorithm of the proposed framework.

A.16 Future Directions

1448

1449

1450

1451

1452

1453 1454

1455

1456

1457

1458

1459

1460

1461

1462

1463

1464This work discovered the possibility of LLMs to1465propose 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).

1466

1467

1468

1469

1470

1471

1472

1473

1474

1475

1476

1477

1478

1479

1480

1481

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.17 Dataset Split of TOMATO

The full dataset is used only as test set.

The license is CC-BY 4.0.

A.19 Method for Prevention of Personal Information

During the dataset collection process, we make1483sure that the dataset is constructed only with public1484information (published papers, Wikipedia, business1485review, and news).1486

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 = Background_Finder(c)
 3:
       if b == None then
         continue
 4:
 5:
       end if
       for iteration k \in 0...M do
 6:
 7:
         if k ! = 0 then
            past_f = Inspiration_Feedback(t, h,
 8:
            present_f)
 9:
         else
            past_f = None
10:
         end if
11:
12:
         t,
                           t\_reason
                                                   =
         Inspiration_Title_Finder(C, b, b_reason,
         past f)
         p = \text{find}_\text{passage}_\text{by}_\text{title}(t, C)
13:
         i = \text{Inspiration}_Finder(b, b_reason, p,
14:
         t_reason)
         s = Hypothesis_Suggestor(b, i)
15:
         h = \text{Hypothesis}_{\text{Proposer}}(b, i, s)
16:
         for iteration t \in 0...N do
17:
            cf, rf, nf = \text{Clarity\_Checker}(h),
18:
            Reality_Checker(h),
            Novelty_Checker(h, S)
19:
            present_f = [cf, rf, nf]
            h = \text{Hypothesis}_{\text{Proposer}}(b, i, s, h,
20:
            present f)
         end for
21:
          H.append(h)
22:
23:
       end for
24: end for
```

```
25: return H
```