
The Meta-RCT Approach to Measuring AI's Labor Market Impact*

Anonymous Author(s)

Affiliation

Address

email

ABSTRACT

1
2 How exposed are different occupations to advances in artificial intelligence (AI)? Existing approaches typically infer exposure from patents,
3 expert surveys, or static mappings between AI capabilities and tasks, but each has limitations. I introduce a complementary framework
4 that aggregates evidence from randomized controlled trials (RCTs) in economics and the social sciences, which directly test the causal
5 impact of AI tools on worker performance. Drawing on RCTs across domains such as writing, coding, forecasting, and management, I map
6 their results to O*NET tasks and then aggregate to the occupational level. The outcome is the 'RCT Exposure Index', a novel, empirically
7 grounded measure of AI exposure that captures heterogeneity across tasks and occupations. By anchoring exposure in experimental
8 evidence rather than proxies, the RCT Exposure Index provides a real-time lens on how AI is reshaping the landscape of work.

Submitted to 1st Open Conference on AI Agents for Science (agents4science 2025). Do not distribute.

*Kris Gulati thanks the Institute of Humane Studies under grant numbers: IHS018498, IHS018315, IHS01854, and Emergent Ventures. We thank Raymond Kim for his advice and feedback. All errors are our own.

11 **KEYWORDS:** AI Benchmarks, Occupational Tasks, Automation Risk, Future of Work, Labor Market Exposure,
12 Task-based Frameworks, Economic Impact of AI

13 *“It is surely a great criticism of our profession that we have not organised a critical summary,*
14 *by specialty or subspecialty, adapted periodically, of all relevant randomised controlled trials.”*
15 – Cochrane (1979)

16 **1 Introduction**

17 How exposed are different occupations to advances in artificial intelligence (AI)? This question has become
18 central to debates over the future of work, productivity, and inequality. A growing body of research has sought
19 to measure occupational exposure to AI and related technologies, but existing approaches face important limi-
20 tations. Measures based on patents Webb (2019) are forward-looking but often sparse, especially since frontier
21 AI labs disclose little intellectual property. Expert surveys Frey and Osborne (2017); Grace et al. (2018) provide
22 valuable intuition but are subjective and infrequently updated. Other approaches use snapshots of AI perfor-
23 mance Felten et al. (2021); Eloundou et al. (2024), which offer a more direct view of technological progress but
24 remain static and disconnected from how AI affects actual worker performance. More recently, Handa et al.
25 (2025) take a usage-based perspective by analyzing millions of real-world Claude conversations mapped to
26 O*NET tasks. Their framework highlights how AI is currently diffusing into work, with especially high penetration
27 in software development and writing but little in roles requiring physical manipulation. In contrast, Gulati (2025)
28 introduce the Benchmark-based AI Occupational Exposure (BAIOE) index, which systematically links progress
29 on frontier AI benchmarks to O*NET tasks and occupations. While usage-based measures emphasize present
30 adoption and benchmark-based measures highlight frontier potential, both remain indirect in capturing how AI
31 actually changes worker productivity.

32 In this paper, I introduce a complementary framework that anchors measurement of AI exposure in experimental
33 evidence. Specifically, I aggregate findings from randomized controlled trials (RCTs) in economics and the
34 social sciences that directly test the causal impact of AI tools on worker performance. These RCTs span
35 diverse domains-including writing, coding, forecasting, and management-providing high-quality evidence on
36 how AI changes speed, accuracy, quality, and decision-making in tasks central to many occupations.

37 To translate these results into labor market relevance, I map each RCT’s estimated treatment effects to O*NET
38 tasks, the standardized taxonomy of work activities and abilities. I then aggregate exposure across tasks to the
39 occupational level, weighting by task importance within each occupation. The outcome is the RCT Exposure
40 Index, a novel, empirically grounded measure of AI exposure that captures heterogeneity across tasks and
41 occupations. Unlike prior measures, the RCT Exposure Index offers a real-time lens on how AI is reshaping
42 work, updating dynamically as new RCTs are conducted.

43 This paper makes three contributions. First, to the best of my knowledge, it is the first to introduce RCTs as
44 a data source for measuring technological exposure, bringing causal identification into a literature often reliant
45 on proxies. Second, it develops a transparent pipeline for linking experimental findings to the O*NET task
46 framework and aggregating them into occupational indices. Third, it demonstrates that the RCT Exposure Index
47 produces meaningful variation across occupations that both complements and challenges existing exposure
48 measures.

49 By anchoring occupational exposure in experimental evidence, this framework highlights where AI is already
50 changing the landscape of work and where gaps remain in our understanding. The RCT Exposure Index is not

51 a substitute for other measures, but a complementary tool that can guide policymakers, researchers, and firms
52 in anticipating how AI diffusion will shape skills, jobs, and economic outcomes.

53 **2 Literature Review**

54 Scholars have long sought to measure how advances in technology affect labor markets. One prominent
55 approach infers exposure from patents and innovation outputs. Webb (2019) uses patent text to measure the
56 overlap between new inventions and occupational descriptions, arguing that patents provide a forward-looking
57 proxy for displacement risk. While effective for many general-purpose technologies, this method has limited
58 traction in AI, where many transformative advances originate in private labs that patent sparingly. Moreover,
59 patents are inherently lagged indicators, reflecting inventions only after formal filings and approvals, often well
60 after the underlying technical progress has occurred.

61 A second line of work relies on expert surveys. Frey and Osborne (2017) famously elicited expert judgments
62 about the susceptibility of occupations to automation, while Grace et al. (2018) surveyed AI researchers to
63 forecast timelines for human-level performance across a wide set of tasks. These studies generate valuable
64 expectations but are inherently subjective, costly, and updated only intermittently.

65 A third, more recent stream grounds measurement in O*NET's task taxonomy. Felten et al. (2021) construct an
66 "AI Occupational Exposure" index by linking expert assessments of AI capabilities to O*NET tasks, while Eloun-
67 dou et al. (2024) extend this approach by connecting the abilities of large language models (LLMs) to O*NET
68 tasks, showing that knowledge-intensive occupations are particularly exposed. These studies move beyond
69 broad occupation-level forecasts by focusing on tasks and underscore the value of transparent frameworks for
70 translating technological performance into economic exposure.

71 Benchmarks have emerged as a complementary tool. Gulati (2025) formalize this approach in the Benchmark-
72 based AI Occupational Exposure (BAIOE) index, directly mapping benchmark performance to O*NET abilities
73 and tasks. Benchmarks provide dynamic, transparent, and continuously updated signals of technical progress,
74 but they remain indirect measures of economic impact since they do not reveal how AI actually affects worker
75 productivity.

76 Another recent direction emphasizes usage. Handa et al. (2025) analyze millions of real-world Claude conver-
77 sations mapped to O*NET tasks, providing an empirical snapshot of how AI tools are being deployed across oc-
78 cupations. Their evidence highlights present-day diffusion -especially in writing and coding - but is constrained
79 by data access and reflects current adoption rather than potential capability. Moreover, usage-based measures
80 face three additional limitations compared to RCTs. First, they capture correlation rather than causal effects:
81 the presence of AI usage in a task does not reveal whether performance improved, worsened, or remained un-
82 changed. Second, usage data may be biased toward early adopters and overrepresented occupations, leading
83 to a skewed picture of diffusion. Third, these measures are often opaque and difficult to replicate, since access
84 to proprietary model logs is restricted. By contrast, RCTs provide transparent, systematically reported, and
85 causally identified evidence of how AI tools change worker performance across tasks.

86 Taken together, these literatures reveal a trade-off. Patents and surveys provide forward-looking but lagged
87 or subjective proxies. Benchmark- and usage-based approaches capture frontier potential or current adoption
88 but remain indirect with respect to worker outcomes. By contrast, randomized controlled trials (RCTs) offer a

89 new, underutilized instrument: they provide causal, transparent, and replicable evidence of how AI alters hu-
 90 man performance on specific tasks. This paper builds on and extends the prior approaches by introducing the
 91 *RCT Exposure Index*, which systematically maps treatment effects from experimental studies to O*NET tasks
 92 and then aggregates them into occupational-level exposure measures. In doing so, it embeds causal identifi-
 93 cation into the measurement of AI's labor market impact, offering a complementary and empirically grounded
 94 alternative to existing proxies.

Study	Data/Approach	Mapping to Work
Webb (2019)	Patent text	Overlap between patents and O*NET tasks
Frey & Osborne (2017)	Expert forecasts	Judgments on occupations' automation risk
Grace et al. (2018)	Expert forecasts	Researcher predictions of AI reaching human-level tasks
Felten et al. (2021)	EFF AI Progress project	10 AI progress areas → O*NET tasks
Eloundou et al. (2024)	LLM evaluations	LLM abilities → O*NET tasks
Tolan et al. (2021/2024)	Research intensity in AI fields	Research intensity in AI fields → Cognitive categories → O*NET tasks
Handa et al. (2025)	Real-world usage data (Claude conversations)	AI usage patterns → O*NET tasks
Gulati (2025)	AI benchmark results	Benchmark performance → O*NET abilities → O*NET tasks
This paper	Randomized controlled trials (RCTs)	Experimental treatment effects on productivity/accuracy → O*NET tasks → Occupations

Table 1: Overview of prior approaches to measuring AI exposure and this paper's contribution.

95 This paper adds to this literature by proposing a new measure for occupational exposure to AI, which we
 96 term the RCT Exposure Index. We build on the growing body of randomized controlled trials that directly
 97 evaluate how AI tools affect human performance across domains such as writing, coding, forecasting, and
 98 management. These studies provide causal evidence on productivity, accuracy, and decision-making, offering
 99 a more rigorous foundation than proxies based on patents, surveys, or benchmarks alone. Because new
 100 RCTs are regularly conducted and reported, their results can be continuously integrated into the framework,
 101 ensuring that measures of exposure remain up to date and empirically grounded. Ultimately, we believe that
 102 leveraging experimental evidence offers the strongest available basis for assessing how AI is reshaping work,
 103 providing a transparent and replicable approach for understanding which occupations are most affected and
 104 how technological diffusion may evolve over time.

105 3 Methodology

106 Our methodological approach of mapping AI benchmark performance into occupational exposure builds on two
 107 recent papers.

108 First, Felten et al. (2021) map the Electronic Frontier Foundation (EFF) AI Progress Measurement project to
 109 ONET tasks through a structured labeling approach. They begin by selecting ten AI applications tracked by
 110 the EFF - such as image recognition, reading comprehension, language modeling, and speech recognition -

111 and then link each application to the 52 workplace abilities defined in O*NET. To establish these links, they
112 run a large survey on Amazon Mechanical Turk, asking respondents (gig workers) to rate how related each AI
113 application is to each O*NET ability. This produces an application-ability relatedness score between 0 and 1
114 for every pairing. The scores are organized into a matrix that systematically connects the EFF applications to
115 O*NET abilities. These ability-level exposures are then aggregated to the occupational level by weighting them
116 with O*NET's measures of ability importance and prevalence, resulting in the AI Occupational Exposure (AIOE)
117 index. This labeling procedure provides a transparent way of mapping progress in frontier AI applications to
118 specific occupational tasks and abilities

119 Second, Eloundou et al. (2024) map large language model (LLM) capabilities to O*NET tasks using a structured
120 labeling approach built around a new "exposure rubric." They begin with O*NET's Detailed Work Activities
121 (DWAs) and tasks, and then assess whether access to an LLM (e.g., via ChatGPT) or to LLM-powered software
122 could reduce the time needed to complete a task by at least 50% without loss of quality. Each task is assigned
123 one of three exposure categories: E0 (no exposure), where LLMs provide minimal benefit or degrade quality;
124 E1 (direct exposure), where the LLM alone can reduce task time by half; and E2 (LLM+ exposure), where the
125 LLM by itself is insufficient, but complementary software or tools built on top of it could achieve that threshold.
126 This study used human annotators to apply the exposure rubric to O*NET tasks, but AI annotators (GPT-4)
127 were found to perform just as well, producing comparable task-level classifications.

128 Building on prior approaches, my framework anchors AI exposure measurement in experimental evidence rather
129 than proxies. Instead of relying on patents, expert surveys, or benchmark scoreboards, I aggregate results from
130 randomized controlled trials (RCTs) in economics and the social sciences that directly test AI's causal effects
131 on human performance. These studies - spanning writing, coding, forecasting, and management - provide
132 treatment effects that are systematically mapped onto O*NET tasks through a structured annotation process.
133 Each RCT result is linked to O*NET abilities, rated for how well the experimental task translates into real-world
134 occupational activities, and scaled by AI's observed impact relative to human performance. These task-level
135 mappings are then aggregated to the occupation level, weighted by O*NET's measures of task importance and
136 prevalence, to produce the RCT Exposure Index.

137 3.1 AI Annotation

138 Our methodology proceeds in a structured sequence designed to translate experimental results from RCTs into
139 task-level measures of occupational exposure. The process consists of four main steps:

140 **Step 1: RCT Identification.** We collected randomized controlled trials in economics and the social sciences
141 that directly test the effects of AI tools on human work. Eligible studies required a clear treatment-control design
142 and reported measurable outcomes related to task performance, such as productivity, accuracy, or decision-
143 making.

144 **Step 2: RCT Summarization.** For each study, the LLM compiled a concise description of the task environ-
145 ment, participant sample, AI system tested, and the reported treatment effect. This ensured a transparent and
146 standardized record of experimental evidence.

147 **Step 3: Task Relevance Mapping.** Each RCT result was mapped to O*NET abilities and tasks using a
148 structured annotation process. Relevance was scored on a 1-3 scale, where 1 = moderate relevance, 2 =

149 strong relevance, and 3 = essential relevance. This rating captures how central the tested ability is to the
150 experimental task and its alignment with real-world occupational activities. This relied on AI annotations.

151 **Step 4: Magnitude of AI's Impact.** For each relevant task, the LLM assigned an effect-size rating based
152 on the study's reported outcomes. Scores ranged from 0 (no effect or negative impact) to 3 (large improve-
153 ment approximately 50% or more), with intermediate categories for small (approximately 10%) and moderate
154 (approximately 25%) improvements. This scale translates heterogeneous outcome measures into a common
155 frame, capturing the causal magnitude of AI's contribution to human performance.

156 Together, these steps provide a systematic method for converting diverse experimental findings into standard-
157 ized task-level annotations, which can then be aggregated to the occupation level.

158 3.2 Aggregation

159 To construct our measure of AI exposure, we integrate two complementary task-level scores into a single
160 index. We begin by assessing task relevance (step 3), which rates how central each O*NET ability is to the
161 experimental task. Next, we evaluate the magnitude of AI's impact (step 4), which captures the causal effect of
162 AI tools on performance for each relevant task. Both scores are normalized onto a common $[0, 1]$ scale and
163 then combined through a weighted average, yielding a composite exposure score at the task level.

164 Following Felten et al. (2021), we then map task-level exposure into occupational exposure. Each task is
165 weighted by O*NET's measures of importance and prevalence, ensuring that frequently performed and central
166 activities exert greater influence on an occupation's overall score. Aggregating in this way links experimental
167 evidence on AI's effects to the heterogeneous distribution of exposure across the labor market.

- 168 • **Relevance score** $R_t \in \{1, 2, 3\}$, where higher values indicate stronger alignment between the ex-
169 perimental task and the O*NET ability.
- 170 • **Impact score** $I_t \in \{0, 1, 2, 3\}$, where higher values indicate larger improvements in human perfor-
171 mance from AI.
- 172 • **O*NET importance** $w_t \in [0, 100]$, the importance of the task within the occupation.
- 173 • **O*NET prevalence** $p_t \in [0, 100]$, the frequency with which the task is performed in the occupation.

174 Task-level exposure

$$E_t = R_t \cdot I_t \cdot w_t \cdot p_t.$$

175 Occupation-level exposure

176 For an occupation consisting of a set of tasks \mathcal{T}_{occ} , the overall exposure is given by:

$$E_{occ} = \sum_{t \in \mathcal{T}_{occ}} E_t.$$

177 4 Results

178 The results presented in Tables 2 and 3 reveal a clear divide in occupational exposure to AI. Table 2 highlights
179 the twenty occupations with the lowest exposure scores. These include dishwashers, cleaners, food service

180 workers, and other manual or service-oriented roles. The low exposure of these occupations aligns with ex-
 181 pectations, since their core activities depend on physical, manual, or interpersonal interaction tasks-domains
 182 where current AI tools have little measurable impact in experimental settings.

183 By contrast, Table 3 lists the twenty occupations with the highest exposure scores, which include scientists
 184 (e.g., physicists, biologists, bioinformaticians), physicians (e.g., emergency medicine, neurology, radiology, gy-
 185 necology), and high-level decision-makers (e.g., chief executives, education administrators). These roles are
 186 knowledge-intensive and rely heavily on analysis, diagnosis, reasoning, and decision-making-areas where ran-
 187 domized controlled trials consistently show substantial productivity improvements from AI.

Table 2: Top 20 Least Exposed Jobs

Rank	Least Exposed
1	Graders and Sorters, Agricultural Products
2	Pressers, Textile, Garment, and Related Materials
3	Dishwashers
4	Slaughterers and Meat Packers
5	Models
6	Cleaners of Vehicles and Equipment
7	Fast Food and Counter Workers
8	Packers and Packagers, Hand
9	Dining Room and Cafeteria Attendants and Bartender Helpers
10	Food Preparation Workers
11	Helpers–Painters, Paperhangers, Plasterers, and Stucco Masons
12	Maids and Housekeeping Cleaners
13	Refuse and Recyclable Material Collectors
14	Shampooers
15	Manicurists and Pedicurists
16	Landscaping and Groundskeeping Workers
17	Locker Room, Coatroom, and Dressing Room Attendants
18	Tapers
19	Foundry Mold and Coremakers
20	Helpers–Brickmasons, Blockmasons, Stonemasons, and Tile and Marble Setters

Notes: Rankings are ordered from 1 (least exposed overall) to 20 (relatively more exposed within this group).

These occupations are identified as having the lowest exposure scores in the dataset.

Table 3: Top 20 Most Exposed Jobs

Rank	Most Exposed
1	Physicists
2	Emergency Medicine Physicians
3	Molecular and Cellular Biologists
4	Biochemists and Biophysicists
5	Preventive Medicine Physicians
6	Neurologists
7	Epidemiologists
8	Physicians, Pathologists
9	Geneticists
10	Neuropsychologists
11	Mathematicians
12	Microbiologists
13	Urologists
14	Bioinformatics Scientists
15	General Internal Medicine Physicians
16	Education Administrators, Kindergarten through Secondary
17	Chief Executives
18	Radiologists
19	Robotics Engineers
20	Obstetricians and Gynecologists

Rankings are ordered from 1 (most exposed overall) to 20 (relatively less exposed within this group).

These occupations are identified as having the highest exposure scores in the dataset.

188 5 Limitations

189 Like any new measure, the RCT Exposure Index faces limitations. First, the availability of RCTs constrains
 190 coverage. At present, most experimental studies of AI focus on knowledge-intensive tasks-writing, coding,
 191 forecasting, and decision-making-while dexterous or manual occupations remain underrepresented. However,
 192 this concentration is not entirely a weakness: by nature, researchers tend to conduct RCTs in domains where
 193 the impact of AI is most likely to be present, where tools are most advanced, and where stakes for performance
 194 are high. As such, the available evidence is still highly informative about the frontier of AI's labor market effects.
 195 Moreover, the framework is designed to update dynamically as new RCTs emerge, progressively expanding
 196 coverage beyond current domains.

197 Second, RCTs vary in design, outcomes, and reporting standards. Some studies emphasize speed, others
 198 quality or accuracy, and still others subjective assessments. While this heterogeneity complicates aggregation,
 199 it is also a strength: the framework is flexible enough to capture multiple dimensions of performance rather than
 200 reducing exposure to a single metric. By systematically standardizing effect sizes and mapping them to O*NET
 201 tasks, the approach ensures transparency and comparability across studies.

202 Third, mapping RCT outcomes to O*NET tasks involves judgment, and some subjectivity is unavoidable. To
 203 mitigate this, the pipeline relies on the standardized O*NET taxonomy and explicit coding procedures that can
 204 be replicated or extended by future researchers. Importantly, the exercise is not meant to claim a single "true"

205 mapping, but to establish a transparent and improvable link between experimental findings and occupational
206 exposure.

207 Finally, while the RCT Exposure Index complements existing approaches, it does not aim to replace them.
208 Patents, expert surveys, and benchmark-based methods capture forward-looking or technical aspects of AI
209 progress that RCTs may miss. The contribution here is to bring experimental evidence-arguably the strongest
210 form of causal inference-into the measurement toolkit. Taken together, these methods provide a more complete
211 and nuanced picture of how AI is reshaping work.

212 **6 Conclusion**

213 This paper has introduced a novel framework for linking experimental evidence from randomized controlled
214 trials (RCTs) to occupational tasks, repositioning RCT findings from isolated academic studies into systematic
215 indicators of AI's labor market impact. By mapping treatment effects onto O*NET abilities, we provide a dynamic
216 method for assessing how AI causally alters human performance at the task and occupation level.

217 Building on prior work, we argue that RCTs represent the strongest available instrument for understanding AI's
218 impact on labor because they identify causal effects rather than relying on proxies that are lagged, subjective,
219 or indirect. Whereas patents, surveys, and benchmarks each capture valuable but partial perspectives, RCTs
220 offer transparent, replicable, and continuously expanding evidence of how AI is reshaping productivity, accuracy,
221 and decision-making. In this sense, our approach extends and strengthens the literature by embedding causal
222 identification into the measurement of occupational exposure.

223 By reframing experimental results as inputs into labor market analysis, this study contributes to a deeper un-
224 derstanding of the future of work and provides a new evidence-based lens for guiding research, policy, and
225 organizational strategy. As AI systems continue to diffuse across domains, our framework offers a systematic
226 foundation for forecasting where automation pressures are most likely to emerge, where human skills remain
227 resilient, and how institutions can design responses that harness innovation while safeguarding workers.

228 **References**

- 229 Cochrane, A. (1979). *Effectiveness and Efficiency: Random Reflections on Health Services*. London: Nuffield
230 Provincial Hospitals Trust.
- 231 Eloundou, T., S. Manning, P. Mishkin, and D. Rock (2024). Gpts are gpts: Labor market impact potential of llms.
232 *Science* 384(6702), 1306–1308.
- 233 Felten, E., M. Raj, and R. Seamans (2021). Occupational, industry, and geographic exposure to artificial
234 intelligence: A novel dataset and its potential uses. *Strategic Management Journal* 42(12), 2195–2217.
- 235 Frey, C. B. and M. A. Osborne (2017). The future of employment: How susceptible are jobs to computerisation?
236 *Technological forecasting and social change* 114, 254–280.
- 237 Grace, K., J. Salvatier, A. Dafoe, B. Zhang, and O. Evans (2018). When will ai exceed human performance?
238 evidence from ai experts. *Journal of Artificial Intelligence Research* 62, 729–754.
- 239 Gulati, K. (2025). Benchmarking the future of work: Mapping ai progress to occupational tasks. *Available at*
240 *SSRN 5452354*.
- 241 Handa, K., A. Tamkin, M. McCain, S. Huang, E. Durmus, S. Heck, J. Mueller, J. Hong, S. Ritchie, T. Belonax,
242 K. K. Troy, D. Amodei, J. Kaplan, J. Clark, and D. Ganguli (2025). Which economic tasks are performed with
243 ai? evidence from millions of claude conversations. *arXiv preprint arXiv:2503.04761*.
- 244 Webb, M. (2019). The impact of artificial intelligence on the labor market. *Available at SSRN 3482150*.

245 **.1 Responsible and Reproducibility Statement**

246 The broader impact of the project is to try and improve the way we can predict technological impact on the labour
247 market. I don't think there are safety considerations involved in this project given the nature of the research.

248 This is a very preliminary project intended for this conference only. If this project is developed, I will work on
249 ensuring that it is well-documented and reproducible, with the code being made publicly available.

250 article

agents4science₂025

251 [utf8]inputenc [T1]fontenc hyperref url booktabs amsfonts nicefrac microtype xcolor

252 **Agents4Science AI Involvement Checklist**

253 1. **Hypothesis development:** Hypothesis development includes the process by which you came to ex-
254 plore this research topic and research question. This can involve the background research performed
255 by either researchers or by AI. This can also involve whether the idea was proposed by researchers
256 or by AI.

257 Answer: B

258 Explanation: The idea was mine (human) and hypothesis development was guided by me but I let the
259 AI do most of the work beyond generating the initial research question

260 2. **Experimental design and implementation:** This category includes design of experiments that are
261 used to test the hypotheses, coding and implementation of computational methods, and the execution
262 of these experiments.

263 Answer: B

264 Explanation: I oversaw most of the methods and empirical exercises, but the AI did a lot of proposing
265 the methods and executing all the labelling etc.

266 3. **Analysis of data and interpretation of results:** This category encompasses any process to organize
267 and process data for the experiments in the paper. It also includes interpretations of the results of the
268 study.

269 Answer: A

270 Explanation: The AI did everything here, i.e. analyses and interpreted the data and results

271 4. **Writing:** This includes any processes for compiling results, methods, etc. into the final paper form.
272 This can involve not only writing of the main text but also figure-making, improving layout of the
273 manuscript, and formulation of narrative.

274 Answer: B

275 Explanation: I oversaw the AI and prompted it in certain directions but almost all of the text is generated
276 by the AI.

277 5. **Observed AI Limitations:** What limitations have you found when using AI as a partner or lead author?

278 Description: Requires human oversight, misses some literature.

279 **Agents4Science Paper Checklist**

280 **1. Claims**

281 Question: Do the main claims made in the abstract and introduction accurately reflect the paper's
282 contributions and scope?

283 Answer: Yes

284 Justification: The abstract and introduction summarise the core method

285 Guidelines:

- 286 • The answer NA means that the abstract and introduction do not include the claims made in the
287 paper.
- 288 • The abstract and/or introduction should clearly state the claims made, including the contributions
289 made in the paper and important assumptions and limitations. A No or NA answer to this question
290 will not be perceived well by the reviewers.
- 291 • The claims made should match theoretical and experimental results, and reflect how much the
292 results can be expected to generalize to other settings.
- 293 • It is fine to include aspirational goals as motivation as long as it is clear that these goals are not
294 attained by the paper.

295 **2. Limitations**

296 Question: Does the paper discuss the limitations of the work performed by the authors?

297 Answer: Yes

298 Justification: We created a limitations section which I think is reasonable

299 Guidelines:

- 300 • The answer NA means that the paper has no limitation while the answer No means that the paper
301 has limitations, but those are not discussed in the paper.
- 302 • The authors are encouraged to create a separate "Limitations" section in their paper.
- 303 • The paper should point out any strong assumptions and how robust the results are to viola-
304 tions of these assumptions (e.g., independence assumptions, noiseless settings, model well-
305 specification, asymptotic approximations only holding locally). The authors should reflect on how
306 these assumptions might be violated in practice and what the implications would be.
- 307 • The authors should reflect on the scope of the claims made, e.g., if the approach was only tested
308 on a few datasets or with a few runs. In general, empirical results often depend on implicit
309 assumptions, which should be articulated.
- 310 • The authors should reflect on the factors that influence the performance of the approach. For
311 example, a facial recognition algorithm may perform poorly when image resolution is low or
312 images are taken in low lighting.
- 313 • The authors should discuss the computational efficiency of the proposed algorithms and how
314 they scale with dataset size.
- 315 • If applicable, the authors should discuss possible limitations of their approach to address prob-
316 lems of privacy and fairness.

- 317
- While the authors might fear that complete honesty about limitations might be used by reviewers as grounds for rejection, a worse outcome might be that reviewers discover limitations that aren't acknowledged in the paper. Reviewers will be specifically instructed to not penalize honesty concerning limitations.
- 318
- 319
- 320

321 **3. Theory assumptions and proofs**

322 Question: For each theoretical result, does the paper provide the full set of assumptions and a complete (and correct) proof?

323

324 Answer: NA

325 Justification: there aren't theoretical assumptions/proofs

326 Guidelines:

- The answer NA means that the paper does not include theoretical results.
 - All the theorems, formulas, and proofs in the paper should be numbered and cross-referenced.
 - All assumptions should be clearly stated or referenced in the statement of any theorems.
 - The proofs can either appear in the main paper or the supplemental material, but if they appear in the supplemental material, the authors are encouraged to provide a short proof sketch to provide intuition.
- 327
- 328
- 329
- 330
- 331
- 332

333 **4. Experimental result reproducibility**

334 Question: Does the paper fully disclose all the information needed to reproduce the main experimental results of the paper to the extent that it affects the main claims and/or conclusions of the paper (regardless of whether the code and data are provided or not)?

335

336

337 Answer: No

338 Justification: I didn't provide all the data and prompting for now because I would like to keep this private to develop a whole paper later.

339

340 Guidelines:

- The answer NA means that the paper does not include experiments.
 - If the paper includes experiments, a No answer to this question will not be perceived well by the reviewers: Making the paper reproducible is important.
 - If the contribution is a dataset and/or model, the authors should describe the steps taken to make their results reproducible or verifiable.
 - We recognize that reproducibility may be tricky in some cases, in which case authors are welcome to describe the particular way they provide for reproducibility. In the case of closed-source models, it may be that access to the model is limited in some way (e.g., to registered users), but it should be possible for other researchers to have some path to reproducing or verifying the results.
- 341
- 342
- 343
- 344
- 345
- 346
- 347
- 348
- 349
- 350

351 **5. Open access to data and code**

352 Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material?

353

354 Answer: No

355 Justification: No, the paper is very preliminary and so the code and data have not been shared,
356 although the method is relatively clear

357 Guidelines:

- 358 • The answer NA means that paper does not include experiments requiring code.
- 359 • Please see the Agents4Science code and data submission guidelines on the conference website
360 for more details.
- 361 • While we encourage the release of code and data, we understand that this might not be possible,
362 so "No" is an acceptable answer. Papers cannot be rejected simply for not including code, unless
363 this is central to the contribution (e.g., for a new open-source benchmark).
- 364 • The instructions should contain the exact command and environment needed to run to reproduce
365 the results.
- 366 • At submission time, to preserve anonymity, the authors should release anonymized versions (if
367 applicable).

368 6. Experimental setting/details

369 Question: Does the paper specify all the training and test details (e.g., data splits, hyperparameters,
370 how they were chosen, type of optimizer, etc.) necessary to understand the results?

371 Answer: No

372 Justification: Again, the paper is still early and preliminary, so I didn't think it made sense to do this
373 just yet because it required a lot more work

374 Guidelines:

- 375 • The answer NA means that the paper does not include experiments.
- 376 • The experimental setting should be presented in the core of the paper to a level of detail that is
377 necessary to appreciate the results and make sense of them.
- 378 • The full details can be provided either with the code, in appendix, or as supplemental material.

379 7. Experiment statistical significance

380 Question: Does the paper report error bars suitably and correctly defined or other appropriate infor-
381 mation about the statistical significance of the experiments?

382 Answer: NA

383 Justification: There are no statistical tests

384 Guidelines:

- 385 • The answer NA means that the paper does not include experiments.
- 386 • The authors should answer "Yes" if the results are accompanied by error bars, confidence inter-
387 vals, or statistical significance tests, at least for the experiments that support the main claims of
388 the paper.
- 389 • The factors of variability that the error bars are capturing should be clearly stated (for example,
390 train/test split, initialization, or overall run with given experimental conditions).

391 8. Experiments compute resources

392 Question: For each experiment, does the paper provide sufficient information on the computer re-
393 sources (type of compute workers, memory, time of execution) needed to reproduce the experiments?

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
421
422
423
424
425

Answer: No

Justification: Again, this is too preliminary to report this in good faith

Guidelines:

- The answer NA means that the paper does not include experiments.
- The paper should indicate the type of compute workers CPU or GPU, internal cluster, or cloud provider, including relevant memory and storage.
- The paper should provide the amount of compute required for each of the individual experimental runs as well as estimate the total compute.

9. Code of ethics

Question: Does the research conducted in the paper conform, in every respect, with the Agents4Science Code of Ethics (see conference website)?

Answer: Yes

Justification: I abided by the ethics of the conference in good faith and reviewed it and followed it accurately

Guidelines:

- The answer NA means that the authors have not reviewed the Agents4Science Code of Ethics.
- If the authors answer No, they should explain the special circumstances that require a deviation from the Code of Ethics.

10. Broader impacts

Question: Does the paper discuss both potential positive societal impacts and negative societal impacts of the work performed?

Answer: Yes

Justification: This is an economics paper and so this is often built into this type of work

Guidelines:

- The answer NA means that there is no societal impact of the work performed.
- If the authors answer NA or No, they should explain why their work has no societal impact or why the paper does not address societal impact.
- Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations, privacy considerations, and security considerations.
- If there are negative societal impacts, the authors could also discuss possible mitigation strategies.