Can LLMs Propose Instrumental Variables for Causal Reasoning?

Anonymous Author(s)
Affiliation
Address
email

Abstract

In the presence of confounding between an endogenous variable and the outcome, instrumental variables (IVs) are used to isolate causal effects. Identifying valid instruments requires interdisciplinary knowledge and contextual understanding, making it a difficult task. In this paper, we examine whether large language models (LLMs) can assist. We adopt a two-stage evaluation: first, testing whether LLMs recover established instruments from the literature, and second, assessing whether they avoid empirically or theoretically discredited ones. Building on these results, we introduce **IV Co-Scientist**, a multi-agent system that proposes, critiques, and refines IVs, along with a statistical test to contextualize consistency without ground truth. Our results show the potential of LLMs to identify valid IVs from large observational data.

1 Introduction

Understanding the causal effect of a treatment on an outcome is a central question across disci-plines [39]. In economics, for example, we may ask: *Does additional schooling increase earnings?*. Estimating such effects is difficult due to endogeneity, when the treatment is correlated with unob-served factors that also influence the outcome [8, 39], and because treatments may be mismeasured or unobservable. To address this, researchers use instrumental variables (IVs)[44, 11, 52]. A valid IV must causally influence the treatment, not be caused by it, and affect the outcome only through the treatment. When these conditions hold, IVs yield consistent causal estimates, even in the presence of unobserved confounders. Identifying valid instruments is therefore crucial yet challenging, as IV strength and validity determine the reliability of causal estimates in domains such as economics[42, 21] and health sciences [12, 5].

Identifying valid instrumental variables is a challenging task that spans theory, domain expertise, and empirical evidence [24]. Statistically, a valid IV must satisfy relevance (correlation with the endogenous treatment) and exclusion (no direct effect on the outcome except through the treatment), while also being independent of unobserved confounders. Meeting these assumptions requires more than statistical tests—it demands contextual knowledge. Experts draw on institutional details, historical context, policy design, or mechanisms from natural and social sciences to argue for or against validity [14]. IV discovery is thus inherently multidisciplinary and often requires creativity, as plausible sources of exogeneity may not be directly observable in data. Moreover, instruments once accepted are sometimes later discredited [36]. For example, rainfall has been used to study the effect of war on national progress, but later research questioned its exclusion validity due to direct effects on outcomes [43]. Such cases underscore both the fragility and difficulty of instrument identification.

Given these challenges, it is natural to ask whether large language models (LLMs) can aid IV discovery. Trained on vast text corpora across economics, health sciences, law, and history, LLMs encode broad domain knowledge that may help generate and assess candidate instruments [9].

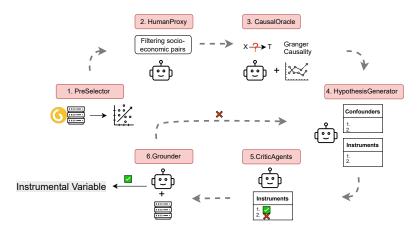


Figure 1: IV Co-Scientist framework, which integrates LLM-based agents with statistical tools.

Recent studies show their utility in literature review, hypothesis generation [45, 2], and experimental design [3, 34], suggesting their promise as tools for early-stage IV discovery. We view LLMs not as substitutes for theoretical reasoning, but as "thinking collaborators" that can extend human intuition and creativity. In this paper, we systematically explore this potential.

LLMs have shown strong performance in causal discovery, often surpassing statistical methods [51, 27, 1]. Their contextual reasoning has been used to identify mediators, extract causal graphs from text, and simulate interventions through structured prompting. We argue that IV discovery is well-suited to LLMs, as it requires both domain-informed reasoning and creative hypothesis generation.

In this paper, we investigate whether LLMs can aid instrumental variable discovery through a structured multi-agent framework, where agents propose, critique, and refine candidate instruments. Our goal is to assess their ability not only to recover established instruments but also to generate novel ones for previously unstudied treatment—outcome pairs. We adopt a staged evaluation: first, testing whether LLMs reproduce well-known instruments from the literature; second, examining whether they avoid those invalidated by theory or evidence; and finally, evaluating their capacity to generate candidate instruments in new contexts, highlighting their potential for IV discovery.

52 2 IV Co-Scientist

Having validated LLMs' ability to recover canonical IVs (subsection Appendix C..1) and avoid discredited ones (subsection Appendix C..2), we next evaluate the system in a fully open-ended setting. Our goal is to test whether LLMs can generate meaningful and potentially novel instrumental variables for real-world causal questions without relying on prior literature. This reflects a realistic and challenging scenario: in applied research, analysts often explore large observational datasets to estimate causal effects for which no established IVs exist, requiring domain expertise, creativity, and data-driven reasoning. We assess whether LLMs, paired with a structured evaluation pipeline, can support this discovery process. Within this pipeline, all validity criteria for IVs are tested: relevance is evaluated statistically, while exclusion and independence are assessed through LLM reasoning, mirroring the approach taken by applied economists. We use a real-world, high-dimensional sandbox to test open-ended causal exploration. The Gapminder dataset [41] includes socio-economic indicators across countries and over time. The dataset has observations for more than 500 such indicators. We aim to find IVs of novel pairs that are still statistically sound.

We formulate a multi-stage, multi-agent system, where each agent is responsible for a specific task in the discovery pipeline (see Figure 2). Let $\mathcal{V} = \{v_1, v_2, \dots, v_n\}$ denote the set of all variables in the dataset. Our goal is to identify a valid IV Z for a treatment-outcome pair $(T, Y) \in \mathcal{V} \times \mathcal{V}$, such that the IV conditions are satisfied. Below, we describe the different stages and agents of IV Co-scientist.

Correlation Filter (*PreSelector*) We compute the Pearson correlation coefficient $\rho(v_i, v_j)$ for all candidate variable pairs, $(v_i, v_j) \in \mathcal{V} \times \mathcal{V}$. We retain pairs satisfying:

$$\mathcal{P} = \{ (v_i, v_j) \mid |\rho(v_i, v_j)| > \tau_{\rho}, \}$$
 (1)

- for pre-defined thresholds τ_{ρ} . This step eliminates weak or statistically noisy pairs. However, since
- 73 correlation strength alone does not account for sample size, we also consider the number of data
- 74 points over which the correlation was computed.
- 75 **Semantic Relevance Agent (HumanProxy)** LLM selects human-meaningful and policy-relevant
- 76 pairs from \mathcal{P} . Output is set $\mathcal{S} = \{(v_i, v_j)\}_s$, to be hypothesized as candidate (T, Y) pairs. This
- step simulates the reasoning a researcher might apply in choosing interpretable or practically socio-
- 78 economic questions.

81

82

83

84

- 79 **Causal Direction Agent** (*CausalOracle*) For each $(v_i, v_j) \in \mathcal{S}$, apply LLM causal reasoning and statistical tests via Granger causality to infer directionality:
 - LLM-based Causal Reasoning: Prompted judgments on whether $v_i \to v_j$ or $v_j \to v_i$, based on world knowledge.
 - Granger Causality Test: Statistical test of the temporal data, verifying whether lagged values of v_i improve the prediction of v_j , beyond v_j 's history. See Appendix F..1 for details.
- We retain only those pairs (v_i, v_j) for which both the LLM and Granger test agree on the direction. The directionally inferred pair is now labeled as (T, Y), with T as treatment and Y as outcome.
- 87 **IV Suggestor Agent** (*HypothesisGenerator*) Given a causal pair $(T \to Y)$, the LLM generates a set 88 of k candidate IVs. See Appendix Appendix C..
- IV Critic Agents (*CriticAgents*) Each candidate Z_i is passed through two critic agents that critique the IVs and give a list of $\mathcal{Z}_{\text{valid}}$. See Appendix Appendix C..

91 Proxy Matching Agent (Grounder)

- For each valid IV $Z_i \in \mathcal{Z}_{\text{valid}}$ proposed by the LLM, we attempt to ground it in the dataset by identifying a concrete proxy variable. If no such proxy is found, Z_i is excluded from downstream
- evaluation. Otherwise, the discovered IV is retained as $(Z_i, \text{Proxy}(Z_i))$.
- 95 The IV Co-scientist starts with the *PreSelector*, which filters variable pairs based on correlation and
- 96 sample size. The *HumanProxy* then selects socio-economically meaningful pairs, forming candidate
- ⁹⁷ causal pairs S. The CausalOracle applies LLM reasoning and Granger tests to infer directionality. For
- 98 each $(T \to Y)$, the *HypothesisGenerator* proposes candidate IVs, which are vetted by *CriticAgents*.
- 99 Valid IVs are then grounded to dataset variables by the *Grounder*, and only those with concrete
- proxies proceed to causal estimation. If no valid IVs were found, then HypothesisGenerator and the
- 101 following agents are rerun.

102 2.1 Evaluation

- Given that the discovered (T, Y, Z) triplets in our open-ended pipeline are novel, direct comparison to
- ground truth IVs is not feasible. To evaluate the plausibility and effectiveness of the LLM-suggested
- 105 IVs, we use the statistical strength of the IV, a standard measure that is used in the IV literature.
- Further, we propose a novel metric to compare sets of valid and invalid IVs.

107 2.1.1 Statistical Strength via F-statistic

- A key requirement for a valid IV is relevance, which means that the IV must be sufficiently correlated
- with the treatment variable. To quantify this, we compute the first-stage F-statistic, a standard method
- used in instrumental variables analysis to detect weak IVs. Specifically, we regress the treatment
- variable T on the candidate IV Z and assess whether Z explains significant variation in T. A high
- F-statistic indicates strong predictive power.
- 113 In our analysis, we use robust heteroskedasticity-consistent estimators that do not assume Gaussian
- errors, reflecting the potentially complex and noisy nature of observational data. The F-statistic tests
- the null hypothesis $H_0: \beta = 0$. An F-statistic value above the conventional threshold (typically 10)
- indicates a strong IV.

117

2.1.2 Consistency of Estimated Effects

- While the relevance of an instrument can be assessed via predictive strength (e.g., F-statistic), its
- overall validity also relies on the more elusive exclusion and independence assumptions, which

	$GDP \rightarrow Health$		Income → Emissions		Sanitation → Mortality		$Poverty \rightarrow Cholesterol$		Female literacy → Kids	
	Relevance	\mathcal{C}_{norm}	Relevance	C_{norm}	Relevance	$C_{ m norm}$	Relevance	C_{norm}	Relevance	$C_{ m norm}$
GPT-4.1	14.28	0.94	17.52	0.91	11.37	0.97	13.44	0.88	19.81	0.93
o3-mini	14.28	0.94	14.88	0.98	11.37	0.97	10.32	0.76	19.81	0.93
QwQ	13.10	0.85	16.05	0.89	10.76	0.94	12.51	0.85	18.23	0.89
Llama3.1 70b	12.65	0.95	15.33	0.89	10.28	0.72	11.90	0.86	17.50	0.79
Llama3.1 8b	13.10	0.85	11.92	0.79	10.76	0.94	12.51	0.85	18.23	0.89

Table 1: Performance of different LLMs in discovering novel IVs. Relevance is defined by the first stage F-statistic and C_{norm} gives the consistency when compared to random IVs.

cannot be directly verified without ground-truth causal effects. To address this, we introduce a novel evaluation metric, **consistency**, to assess the quality of IV sets.

The idea is that if a set of instruments truly isolates exogenous variation in the treatment, each 122 instrument should produce similar estimates of the average treatment effect (ATE) via two-stage 123 least squares (2SLS). For two instruments Z_1 and Z_2 proposed by the LLM, we compute their 124 2SLS estimates $\hat{\beta}_{ATE}^{(Z_1)}$ and $\hat{\beta}_{ATE}^{(Z_2)}$ and define the consistency score as $\Delta_{LLM} = \left|\hat{\beta}_{ATE}^{(Z_1)} - \hat{\beta}_{ATE}^{(Z_2)}\right|$ 125 where smaller values indicate stronger internal agreement. To contextualize this measure, we 126 construct a null distribution by randomly sampling proxy variables R_1 and R_2 from the dataset, 127 defining $\Delta_{\text{Rand}} = \left| \hat{\beta}_{\text{ATE}}^{(R_1)} - \hat{\beta}_{\text{ATE}}^{(R_2)} \right|$, which captures the variability expected from invalid or spurious 128 instruments. Comparing Δ_{LLM} to Δ_{Rand} allows us to test whether LLM-suggested IVs exhibit greater 129 internal consistency than would be expected by chance. Inspired by the self-compatibility test in 130 causal discovery [16], this approach provides indirect evidence that LLMs may identify variables 131 isolating genuine exogenous variation, even when the exclusion and independence assumptions 132 cannot be directly verified. 133

2.2 Results

134

In Table 1, we evaluate the quality of LLM-suggested IVs from two perspectives: statistical relevance and consistency of estimated causal effects. An IV is considered strong if it predicts variation in the treatment T, typically with an F-statistic exceeding 10. We examine five examples autonomously generated by our multi-agent IV Co-Scientist (see Appendix Appendix H. for details). To quantify stability, we define a normalized consistency score $C_{\text{norm}} = \left| \frac{\Delta_{\text{LLM}} - \Delta_{\text{Rand}}}{\Delta_{\text{Rand}}} \right|$, where values near 1 indicate that LLM-suggested IVs are more internally consistent than random proxies.

Many LLM-suggested IVs achieve high relevance, though statistical significance alone does not guarantee sufficient strength to avoid weak-IV issues [29]; hence, we emphasize C_{norm} . Empirically, C_{norm} scores are often near or above 1, showing a clear gap between LLM-suggested and random IVs, with GPT-4.1 and o3-mini producing similar IVs and results. Figure 3 further supports this: panel (a) shows posterior distributions of ATE1 and ATE2 using LLM-suggested IVs, suggesting both relevance and meaningful local treatment effect heterogeneity, whereas panel (b) highlights weak or inconsistent IVs.

Human evaluation. We consulted a faculty-level economist to qualitatively assess the LLMgenerated IVs. They found the *CriticAgents*' reasoning and confounder identification generally sound.
They noted that accepted and rejected IVs often differ not in validity but in generality: accepted ones tend to be broader and less debated, while rejected ones are more specific, often with known critiques.

3 Conclusion

152

153

154

155

156

157

158

Instrumental variables are central to causal inference in observational studies, but identifying them is challenging and typically demands deep domain expertise. While large language models offer new opportunities for extracting knowledge from text, their use in discovering instruments beyond toy examples remains underexplored. We introduce a multi-agent framework that analyses the data, proposes candidate instruments for a given treatment-outcome pair, and validates them semantically. In addition, we propose a consistency-based metric to assess internal validity in the absence of ground truth. Our empirical results on real-world data demonstrate that LLM-suggested instruments show meaningful consistency, providing a first step toward principled use of LLMs in variable discovery.

References

- [1] Ahmed Abdulaal, adamos hadjivasiliou, Nina Montana-Brown, Tiantian He, Ayodeji Ijishakin
 and Ivana Drobnjak, Daniel C. Castro, and Daniel C. Alexander. Causal modelling agents:
 Causal graph discovery through synergising metadata- and data-driven reasoning. In *ICLR*,
 2024.
- [2] Shubham Agarwal, Gaurav Sahu, Abhay Puri, Issam H Laradji, Krishnamurthy Dj Dvijotham,
 Jason Stanley, Laurent Charlin, and Christopher Pal. Litllms, llms for literature review: Are we
 there yet? *Transactions on Machine Learning Research*.
- [3] Microsoft Research AI4Science and Microsoft Azure Quantum. The impact of large language models on scientific discovery: a preliminary study using gpt-4. *arXiv*, 2023.
- [4] Joshua D Angrist and Alan B Krueger. Instrumental variables and the search for identification: From supply and demand to natural experiments. *Journal of Economic perspectives*, 15(4):69–85, 2001.
- [5] Michael Baiocchi, Jing Cheng, and Dylan S Small. Instrumental variable methods for causal inference. *Statistics in medicine*, 33(13):2297–2340, 2014.
- 176 [6] Taiyu Ban, Lyvzhou Chen, Xiangyu Wang, and Huanhuan Chen. From query tools to causal architects: Harnessing large language models for advanced causal discovery from data. *arXiv*, 2023.
- [7] Shraddha Barke, Michael B James, and Nadia Polikarpova. Grounded copilot: How programmers interact with code-generating models. *Proceedings of the ACM on Programming Languages*, 7(OOPSLA1):85–111, 2023.
- [8] Roger John Bowden, Roger J Bowden, and Darrell A Turkington. *Instrumental variables*.
 Number 8. Cambridge university press, 1990.
- [9] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal,
 Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are
 few-shot learners. In *NeurIPS*, 2020.
- [10] S. Burgess, Dylan S. Small, and S. Thompson. A review of instrumental variable estimators for mendelian randomization, 2015.
- 189 [11] Raymond J Carroll and Leonard A Stefanski. Measurement error, instrumental variables and corrections for attenuation with applications to meta-analyses. *Statistics in Medicine*, 13(12):1265–1282, 1994.
- 192 [12] John Cawley and Chad Meyerhoefer. The medical care costs of obesity: an instrumental variables approach. *Journal of health economics*, 31(1):219–230, 2012.
- 194 [13] Victor-Alexandru Darvariu, Stephen Hailes, and Mirco Musolesi. Large language models are effective priors for causal graph discovery. *arXiv*, 2024.
- [14] Neil M Davies, George Davey Smith, Frank Windmeijer, and Richard M Martin. Issues in the
 reporting and conduct of instrumental variable studies: a systematic review. *Epidemiology*,
 24(3):363–369, 2013.
- 199 [15] Richard Disney, John Gathergood, Stephen Machin, and Matteo Sandi. Does homeownership reduce crime? a radical housing reform from the uk. *The Economic Journal*, 133(655):2640–201 2675, 2023.
- Philipp M Faller, Leena C Vankadara, Atalanti A Mastakouri, Francesco Locatello, and Dominik
 Janzing. Self-compatibility: Evaluating causal discovery without ground truth. In *International Conference on Artificial Intelligence and Statistics*, pages 4132–4140. PMLR, 2024.
- 205 [17] Trevor S. Gallen. Broken instruments, 2020.
- 206 [18] Gautam Gowrisankaran and Robert J Town. Estimating the quality of care in hospitals using instrumental variables. *Journal of health economics*, 18(6):747–767, 1999.

- 208 [19] Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, et al. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*, 2024.
- 211 [20] Sukjin Han. Mining causality: Ai-assisted search for instrumental variables. *arXiv preprint* arXiv:2409.14202, 2024.
- 213 [21] James Heckman and Salvador Navarro-Lozano. Using matching, instrumental variables, and control functions to estimate economic choice models. *Review of Economics and statistics*, 86(1):30–57, 2004.
- 216 [22] Binyuan Hui, Jian Yang, Zeyu Cui, Jiaxi Yang, Dayiheng Liu, Lei Zhang, Tianyu Liu, Jia-217 jun Zhang, Bowen Yu, Keming Lu, et al. Qwen2. 5-coder technical report. *arXiv preprint* 218 *arXiv:2409.12186*, 2024.
- [23] Guido W Imbens and Paul R Rosenbaum. Robust, accurate confidence intervals with a weak
 instrument: quarter of birth and education. *Journal of the Royal Statistical Society Series A:* Statistics in Society, 168(1):109–126, 2005.
- 222 [24] Wei Jiang. Have instrumental variables brought us closer to the truth. *Review of Corporate Finance Studies*, 6(2):127–140, 2017.
- [25] Carlos E Jimenez, John Yang, Alexander Wettig, Shunyu Yao, Kexin Pei, Ofir Press, and
 Karthik R Narasimhan. Swe-bench: Can language models resolve real-world github issues? In
 The Twelfth International Conference on Learning Representations.
- 227 [26] Yiqiao Jin, Qinlin Zhao, Yiyang Wang, Hao Chen, Kaijie Zhu, Yijia Xiao, and Jindong Wang.
 228 Agentreview: Exploring peer review dynamics with llm agents. In *Proceedings of the 2024*229 *Conference on Empirical Methods in Natural Language Processing*, pages 1208–1226, 2024.
- 230 [27] Emre Kıcıman, Robert Ness, Amit Sharma, and Chenhao Tan. Causal reasoning and large language models: Opening a new frontier for causality. *arXiv*, 2023.
- 232 [28] Apoorva Lal, Mac Lockhart, Yiqing Xu, and Ziwen Zu. How much should we trust instrumental variable estimates in political science? practical advice based on over 60 replicated studies. *arXiv preprint arXiv:2303.11399*, 2023.
- [29] Daniel J Lewis and Karel Mertens. A robust test for weak instruments with multiple endogenous
 regressors. Technical report, Staff Reports, 2022.
- [30] Long Li, Weiwen Xu, Jiayan Guo, Ruochen Zhao, Xingxuan Li, Yuqian Yuan, Boqiang Zhang,
 Yuming Jiang, Yifei Xin, Ronghao Dang, et al. Chain of ideas: Revolutionizing research via
 novel idea development with llm agents. arXiv preprint arXiv:2410.13185, 2024.
- 240 [31] Ethan Lin, Zhiyuan Peng, and Yi Fang. Evaluating and enhancing large language models for novelty assessment in scholarly publications. In *Annual Conference of the Nations of the* 242 *Americas Chapter of the Association for Computational Linguistics*, page 46, 2025.
- [32] Joan Llull. The effect of immigration on wages: exploiting exogenous variation at the national level. *Journal of Human Resources*, 53(3):608–662, 2018.
- 245 [33] Stephanie Long, Tibor Schuster, Alexandre Piché, ServiceNow Research, et al. Can large language models build causal graphs? arXiv, 2023.
- ²⁴⁷ [34] Chris Lu, Cong Lu, Robert Tjarko Lange, Jakob Foerster, Jeff Clune, and David Ha. The ai scientist: Towards fully automated open-ended scientific discovery. *arXiv*, 2024.
- [35] Pingchuan Ma, Rui Ding, Shuai Wang, Shi Han, and Dongmei Zhang. Demonstration of insight pilot: An Ilm-empowered automated data exploration system. arXiv preprint arXiv:2304.00477,
 2023.
- 252 [36] Jonathan Mellon. Rain, rain, go away: 194 potential exclusion-restriction violations for studies using weather as an instrumental variable. *American Journal of Political Science*, 2024.
- 254 [37] OpenAI. Gpt-4o. 2025.

- 255 [38] OpenAI. o3-mini. 2025.
- 256 [39] Judea Pearl. Causality. Cambridge university press, 2009.
- Linlu Qiu, Liwei Jiang, Ximing Lu, Melanie Sclar, Valentina Pyatkin, Chandra Bhagavatula, Bailin Wang, Yoon Kim, Yejin Choi, Nouha Dziri, et al. Phenomenal yet puzzling: Testing inductive reasoning capabilities of language models with hypothesis refinement. In *ICLR*, 2024.
- 260 [41] Hans Rosling et al. Gapminder, for a fact-based worldview. *Gapminder, Stockholm, Sweden.* [online] URL: http://www.gapminder.org, 2010.
- [42] John D Sargan. The estimation of economic relationships using instrumental variables. *Econometrica: Journal of the econometric society*, pages 393–415, 1958.
- Heather Sarsons. Rainfall and conflict: A cautionary tale. *Journal of development Economics*, 115:62–72, 2015.
- 266 [44] Susanne M Schennach. Recent advances in the measurement error literature. *Annual review of economics*, 8(1):341–377, 2016.
- [45] Dmitry Scherbakov, Nina Hubig, Vinita Jansari, Alexander Bakumenko, and Leslie A Lenert.
 The emergence of large language models (llm) as a tool in literature reviews: an llm automated systematic review. arXiv preprint arXiv:2409.04600, 2024.
- 271 [46] Amit Sharma. Necessary and probably sufficient test for finding valid instrumental variables. 272 arXiv preprint arXiv:1812.01412, 2018.
- [47] Ivaxi Sheth, Sahar Abdelnabi, and Mario Fritz. Hypothesizing missing causal variables with llms. *arXiv preprint arXiv:2409.02604*, 2024.
- [48] Ivaxi Sheth, Bahare Fatemi, and Mario Fritz. Causalgraph2llm: Evaluating llms for causal queries. *arXiv preprint arXiv:2410.15939*, 2024.
- 277 [49] Chenglei Si, Diyi Yang, and Tatsunori Hashimoto. Can llms generate novel research ideas? a large-scale human study with 100+ nlp researchers. *ICLR*, 2025.
- [50] Dylan S Small and Paul R Rosenbaum. War and wages: the strength of instrumental variables
 and their sensitivity to unobserved biases. *Journal of the American Statistical Association*,
 103(483):924–933, 2008.
- Aniket Vashishtha, Abbavaram Gowtham Reddy, Abhinav Kumar, Saketh Bachu, Vineeth N
 Balasubramanian, and Amit Sharma. Causal inference using llm-guided discovery. arXiv, 2023.
- ²⁸⁴ [52] Tom Wansbeek and Erik Meijer. Measurement error and latent variables. *A companion to theoretical econometrics*, pages 162–179, 2001.
- [53] Fangzhi Xu, Qika Lin, Jiawei Han, Tianzhe Zhao, Jun Liu, and Erik Cambria. Are large
 language models really good logical reasoners? a comprehensive evaluation from deductive,
 inductive and abductive views. arXiv, 2023.

289 Appendix A. Related Works

LLMs and Causality. LLMs have been used as priors to discover the relationships between causal variables [27, 33, 13, 6, 51, 1, 48]. These methods, alone or in combination with statistical or deep learning methods outperformed the latter. Sheth et al. [47] proposed a benchmark for discovering causal variables from a partial graph; however, they compared against the ground truth semantically without knowing the statistical effect. The closest work to ours is by Sukjin Han [20], however, our work goes beyond prompting the previously established instruments.

LLMs and Scientific Discovery. LLMs are increasingly integrated into various stages of the scientific research workflow, including hypothesis generation and reasoning [53, 40, 49, 30], coding and implementation [25, 7], data analysis [35], and even peer review [26]. Despite their growing role, it remains challenging to assess the significance or scientific plausibility of hypotheses generated by LLMs [31], especially when another LLM is used as a judge. In this work, we study LLM-driven discovery of instrumental variables for novel treatment—outcome pairs and propose evaluation metrics to validate their statistical and causal validity.

Testing an instrumental variable beyond the relevance test is challenging [46], hence we introduce a treatment effect-based consistency metric to quantify the stability of causal estimates across candidate instruments, offering indirect evidence of validity.

306 Appendix B. Preliminaries

The instrumental variable enables identification of causal effects in the presence of endogeneity, i.e, when the treatment variable is correlated with unobserved confounders that also influence the outcome. Formally, let T denote the treatment variable, Y the outcome, and U represent unobserved confounders. A valid instrument Z is a variable that influences T but does not directly affect Y, except through its effect on T.

We consider the structural equations:

$$T = f(Z, U_T) (2)$$

$$Y = g(T, U_Y) \tag{3}$$

where U_T and U_Y may be arbitrarily dependent due to shared unobserved variables U. In this setup, Z qualifies as an instrumental variable for estimating the causal effect of T on Y if the IV validity conditions [39] are satisfied.

316 Appendix B..1 IV Validity Conditions

Relevance. The instrument must be predictive of the treatment. Formally, Z must have a non-zero association with T: $Cov(Z,T) \neq 0$. The relevance condition implies that the function $f(Z,U_T)$ is not a constant in Z, hence Z has a causal effect on T. In practice, this condition is assessed through the first-stage regression.

Exclusion Restriction. The instrument must not directly affect the outcome, nor through any path other than via T. Formally, this means that Z is conditionally independent of Y given T and any covariates $X: Y \perp \!\!\! \perp Z \mid T, X$. This assumption cannot be empirically tested and is usually justified via domain knowledge.

Independence. The instrument must be conditionally independent of the unobserved confounders: $Z \perp \!\!\! \perp U \mid X$. This ensures that the variation in T induced by Z is as random concerning the potential outcomes. Similar to the exclusion criteria, independence is also usually argued by domain knowledge.

Appendix B..2 Estimation via Two-Stage Least Squares

When a valid instrument is available, the causal effect of T on Y can be estimated consistently using Two-Stage Least Squares (2SLS). Following the economics literature [4], we focus on the linear model. This involves:

• Stage 1: Regress T on Z (and covariates X) to obtain predicted treatment \hat{T} :

$$T = \alpha_0 + \alpha_1 Z + \alpha_2 X + \varepsilon_T$$

• Stage 2: Regress Y on the fitted values \hat{T} (and X):

$$Y = \beta_0 + \beta_1 \hat{T} + \beta_2 X + \varepsilon_Y \text{s.t.}$$
 $\varepsilon_Y \perp (Z, X)$

Under the IV assumptions, the coefficient β_1 consistently estimates the causal effect of T on Y. For the rest of the paper, we assume that the ε we focus on is linear noise.

It is important to note that while the relevance condition is statistically testable, the exclusion and independence conditions are not, and must be argued through theory, domain expertise, or natural experiments. This makes the process of identifying valid IVs fundamentally interdisciplinary and often creative. As such, the search for IVs can benefit from tools that integrate reasoning, background knowledge, and flexible hypothesis generation, a role LLMs may be suited to play.

Appendix C. LLMs for IV Reasoning

Our goal is to explore whether LLMs can assist in the discovery of novel and valid IVs. We propose a multi-agent pipeline that separates the creative and evaluative stages of IV discovery, mirroring how human researchers hypothesize and then vet candidate IVs. Given a treatment–outcome pair (T,Y), we define a two-stage LLM-based framework. In the first stage, HypothesisGenerator, two agents are prompted with a causal query to propose candidate instruments $\{Z_1,\ldots,Z_i\}$ and confounders $\{U_1,\ldots,U_j\}$. In addition to hypothesis generation, it is essential to have LLMs act as proxy domain experts to reason about statistically untestable conditions. Thus, the second stage, CriticAgents, involves two LLMs independently evaluating validity: one assesses the exclusion restriction, i.e., whether Z_i affects Y only through T, and the other assesses independence, i.e., whether Z_i is independent of unobserved confounders U_j that influence both T and Y. Each candidate instrument Z_i receives binary feedback from both agents, and only those satisfying both conditions are retained, i.e., $Z_{\text{valid}} = \{Z_i \mid \text{Ex}(Z_i) \land \text{Ind}(Z_i)\}$.

Model	Military se EM ↑	ervice → Earning CM ↑	Education	n → Wages CM ↑	Housing EM ↑	$CM \uparrow Crime$	Healthcar EM ↑	$ \begin{array}{c} \text{re} \rightarrow \text{Mortality} \\ \text{CM} \uparrow \end{array} $	Migratio EM ↑	$\begin{matrix} \mathbf{n} \to \mathbf{Wages} \\ \mathrm{CM} \uparrow \end{matrix}$
GPT-40	0.74	1.00	0.82	1.00	0.75	0.83	0.68	0.91	0.40	0.74
o3-mini	0.73	1.00	0.82	1.00	0.37	0.53	0.59	0.89	0.45	0.81
QwQ	0.74	1.00	0.73	1.00	0.39	0.75	0.52	0.90	0.31	0.70
Llama3.1 8B	0.28	0.42	0.48	0.76	0.36	0.49	0.32	0.65	0.35	0.60
Llama3.1 70B	0.61	0.84	0.67	1.00	0.59	0.75	0.52	0.83	0.57	0.77

Table 2: Performance of LLMs in recovering canonical instrumental variables across five benchmark treatment-outcome pairs. Exact Match (EM) captures direct or paraphrased mentions of literature-established IVs, while Conceptual Match (CM) identifies plausibly equivalent proxies judged by an LLM critic.

Appendix C..1 Recovering Canonical IVs

It is essential to first assess whether they can recover IVs that are already well-established in the literature before talking about novel instruments. This serves two purposes: (1) it helps calibrate the LLM's alignment with existing scientific knowledge and reasoning, and (2) it provides a baseline for evaluating the model's ability to reason causally and contextually about treatment-outcome relationships. If an LLM is unable to identify canonical instruments, then relying on it for more speculative and novel discovery becomes difficult to justify.

We curate a benchmark dataset consisting of well-studied treatment-outcome pairs from economics, health sciences, and social sciences, where valid instrumental variables have been previously proposed and accepted in the literature. Each entry in the benchmark includes a treatment variable T (for example, years of schooling), an outcome variable Y (such as future earnings), and one or more canonical instrumental variables $\{Z_1^*, Z_2^*, \ldots\}$ sourced from peer-reviewed literature.

Model	GDP — HG↓	→ Conflict Critic ↓	BMI HG↓	→ SBP Critic↓	Church HG↓	→ Crime Critic ↓	Turnout HG↓	→ Vote Share Critic ↓	Protest HG↓	$s \rightarrow Prices$ Critic \downarrow
GPT-40	1	0	1	1	1 1	0	0	0	1	0
o3-mini	1	0	1	0	1	1	1	0	1	1
QwQ	1	0	1	1	1	0	1	1	1	1
Llama3.1 8B	0	1	0	1	1	0	0	0	1	0
Llama3.1 70B	1	0	1	0	1	0	1	1	1	0

Table 3: Performance of different LLMs in identifying flawed instruments across treatment—outcome pairs. HG indicates whether the HypothesisGenerator proposed flawed IV and Critic when CriticAgent successfully picks an invalid IV.

367 Appendix C..2 Avoiding Invalid IVs

While the ability to recover canonical instruments is important, an equally critical aspect of evaluating LLMs for instrumental variable discovery is their sensitivity to invalid instruments. Several variables that were historically proposed as instruments have been subsequently discredited due to theoretical objections or empirical evidence, typically involving violations of the exclusion restriction or the independence assumption. For example, IV like rainfall had been used to estimate of the effect of economic activity on civil conflict, but later critiques have revealed direct causal paths or unmeasured confounding, undermining their validity [36].

In this experiment, we aim to assess whether LLMs are able to avoid suggesting such invalidated instruments when prompted with the original treatment outcome pair. This evaluation probes the depth of the model's reasoning: does it simply retrieve past associations?

We design a multi-stage evaluation framework to assess the robustness of LLMs in handling invalid instruments. Our goal is: (1) to test whether the LLM proposer avoids historically invalid instruments on its own, and (2) to evaluate whether the critic LLM can reliably detect and reject such instruments, even when explicitly introduced.

Given a treatment-outcome pair (T, Y) with a documented invalid instrument Z^- (e.g., rainfall), we perform the following steps:

- 1. **Proposer Behavior.** We prompt the LLM proposer to generate a list of k candidate instruments $\{Z_1, Z_2, \ldots, Z_i\}$. We then evaluate whether the model reproduces Z^- or semantically equivalent variants. This allows us to assess whether the proposer model has internalized the criticisms of certain instruments or simply replicates canonical (yet flawed) examples from the literature.
- 2. **Critic Evaluation.** Regardless of Stage 1, we now explicitly inject Z^- into the list of candidate instruments. This injected list is:

$$\{Z_1, Z_2, Z_3, Z^-, Z_4, Z_5\}$$
 (4)

We pass this set through the *CriticAgents*, each independently evaluating the instrument on the Exclusion and Independence criteria.

Appendix C...3 Results

384

385

386

387

388

389

390

391

392

393

We evaluate a range of benchmark LLMs to assess their ability to propose and critique instrumental 394 variables. For the generation stage, we test both reasoning models: o3-mini [38] and QwQ [22] and 395 standard models: GPT-40 [37], Llama3.1 8B [19], and Llama3.1 70B [19]. The HypothesisGenerator 396 then evaluates each candidate instrumental variable (IV) from the generated list Z_1, \ldots, Z_i along 397 with the CriticAgent validating them. The exclusion check is performed independently for each IV, 398 while the independence check is done via comparisons between each IV and the set of hypothesized 399 confounders. We fix i = j = 5 as a balance between promoting diversity in generation and 400 maintaining computational efficiency. We prompt all models with an economist persona to elicit 401 appropriate reasoning. 402

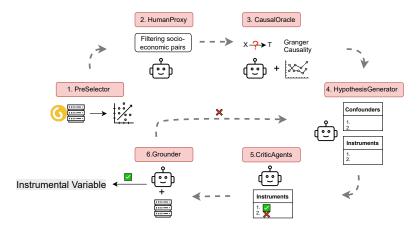


Figure 2: Overview of the *IV Co-Scientist* framework, which integrates LLM-based agents with traditional statistical tools.

403 Appendix C..3.1 Recovering Canonical IVs.

We consider treatment, outcome, and instrument tuples (T,Y,Z^*) sourced from literature. In particular, the outcome of earnings due to military service [50], the effect of education on wages [23], housing and its effect on crime [15], healthcare on mortality [18], and migration's effect on wages [32]. For each pair (T,Y), we prompt the multi-agent HypothesisGenerator and CriticAgents to generate a set of LLM validated candidate instruments \mathcal{Z}_{valid} . We then compare these IVs to the known literature instruments $\{Z_1^*, \ldots\}$ using two matching strategies:

- 1. **Exact Match (EM):** Semantic similarity checks to identify if a known instrument is directly mentioned or closely paraphrased.
 - 2. **Conceptual Match (CM):** LLM-judge whether a generated candidate is a plausible conceptual equivalent or proxy to the known instrument.

Table 2 summarizes the ability of different models to recover well-established instrumental variables across five canonical treatment-outcome settings. We observe that the strongest models: GPT-40, o3-mini, and QwQ2.5 can recover canonical instruments with high consistency. Across all of the settings, we observe CM rating is higher than EM, because while LLMs often propose valid instruments that align with the underlying causal rationale, they frequently use alternate phrasings or suggest closely related proxies.

Appendix C..3.2 Avoiding invalid IV

410

411

412

413

420

421 Given that we have observed positive results in subsubsection Appendix C...3.1, we are interested in 422 evaluating whether LLMs can recognize and avoid historically discredited IVs. We evaluate whether they suggest the negative IV Z^- and whether the *CriticAgents* can filter them out. We filter these 423 (T, Y) from established literature. In particular, the effect of GDP on conflict in a country [36], 424 the effect of body mass index on systolic blood pressure [10], the effect of church attendance on 425 crime [17], the effect of vote turnout on party vote share [28], and protests on consumer prices [36]. 426 In Table 3, we evaluate how well different models handle flawed instruments. The HG column 427 indicates whether the model directly proposed the flawed instrument (Z^-) , while the Critic column 428 captures whether the *CriticAgent* correctly identified and flagged the flaw. 429

Overall, we see that the *CriticAgent* plays a vital role in safeguarding against invalid instruments.

Even when powerful models like GPT-40 and QwQ occasionally suggest flawed variables, the critic is often able to detect and reject them. This highlights the utility of incorporating an automated critic to evaluate statistical validity post hoc. Interestingly, the Llama3.1 8B model appears more conservative, doesn't propose many flawed IVs. However, when such variables are injected, its critic fails to detect the issue.

Appendix D. Experimental Setup

In our paper, we conduct extensive experiments across reasoning and non-reasoning models.

Appendix D..1 EM and CM 438

- In Section Appendix C..1, we evaluate the ability of LLMs to recover canonical instrumental variables using two metrics, building upon [47]: Exact Match (EM) and Conceptual Match (CM).
- **Exact Match (EM).** This metric quantifies the semantic similarity between each LLM-suggested 441 instrument and the known ground-truth instrument. Specifically, we embed both the LLM's suggestion 442
- and the literature-sourced IV using the Qwen3-Embedding-0.6B model. We then compute the cosine similarity between the embeddings and report the similarity score, where higher values indicate closer
- semantic alignment. 445
- Conceptual Match (CM). Exact matches may underestimate the utility of LLM suggestions that 446
- are plausible but lexically dissimilar. To account for this, we introduce a softer, human-grounded 447
- measure: Conceptual Match (CM). For each LLM-generated IV, we prompt another LLM to act as 448
- a domain-aware judge and rate—on a scale from 1 to 10—how conceptually similar the suggestion is 449
- to the accepted IV in terms of causal plausibility. A score closer to 10 indicates a stronger conceptual 450
- match. 451

444

- Together, EM and CM allow us to evaluate both surface-level and deeper, contextual alignment 452
- between LLM-suggested and literature-backed instruments. 453

Appendix D..2 Gapminder Dataset 454

- To evaluate the ability of LLMs to propose and validate novel instrumental variables, we require a 455
- 456 rich and diverse source of real-world observational data. For this purpose, we use the Gapminder
- database¹, a curated compilation of time-series indicators covering over 200 countries and territories. 457
- Gapminder provides over 500 socio-economic and health-related variables, including measures 458
- such as GDP per capita, life expectancy, sanitation access, education levels, and fertility rates. 459
- These indicators are compiled from authoritative sources like the World Bank, WHO, and UN, and 460
- are harmonized to ensure consistency across countries and years. It is under Creative Commons 461
- Attribution 4.0. 462
- The diversity and breadth of variables make Gapminder particularly well-suited for causal analysis. It 463
- contains plausible treatment and outcome variables across multiple domains (e.g., income, health, de-464
- mographics), along with a large pool of potential proxy variables. Moreover, the data are longitudinal, 465
- enabling time-aware causal reasoning techniques such as Granger causality. 466
- We extract country-year level data for all variables with sufficient temporal coverage. To preprocess it, 467
- we removed datapoints with missing values and standardized variable scales. This yields a structured 468
- dataset suitable for evaluating both traditional statistical tests (e.g., relevance via F-statistic) and
- 470 novel LLM-generated instruments under realistic conditions.

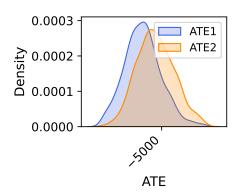
Appendix D..3 Compute 471

- We ran Qwen2.5, LLama70b and Llama 8b on A100 GPUs. The GPT models were accessed via API. 472
- The PreSelector, HumanProxy and CausalOrale just had to run once, while the discovery modules 473
- were iterated over all examples and for new IVs. 474

Appendix D..4 Reproducibility 475

- All LLMs were run with a temperature of 0 and top-p of 1 to ensure deterministic outputs. The results 476
- reported in Table 2 reflect averaged metrics over multiple runs where applicable. Table 3 contains no 477
- variance due to LLM randomness, as the models were only used to suggest instruments, which were 478
- then evaluated against fixed proxies or ground truth. 479

¹https://www.gapminder.org/data/



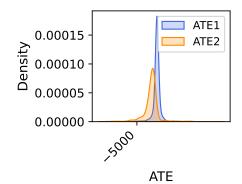


Figure 3: Comparison of the ATE density while using two different IVs: (a) LLM proposed and (b) random. This is for Sanitation \rightarrow Mortality for GPT-40.

We commit to releasing all code, prompts, and evaluation scripts upon acceptance to support full reproducibility.

482 Appendix E. Qualitative Analysis

To complement our quantitative analyses, we conducted a structured qualitative evaluation with an economics professor and political scientist familiar with instrumental variables. The goal was to assess the plausibility and relevance of LLM-generated variables through expert judgment. The expert was given a short document, consisting of four tasks:

Appendix E..1 Task Overview

- Task 1: Generator vs. Critic Evaluation. The expert was asked whether they agreed with the LLM's rejection of certain candidate instruments following critique by a secondary "critic" model that evaluated IVs based on the standard assumptions of relevance, independence, and exclusion.
- Task 2: Agreement with Accepted and Rejected Instruments. The expert reviewed a table of treatment—outcome pairs, each with a list of instruments accepted or rejected by the LLM pipeline. They were asked to indicate agreement with each group of variables (e.g., 2/3 accepted IVs, 1/3 rejected IVs).
- Task 3: Case-Based Evaluation. The expert was presented with a specific example—female literacy as a treatment for fertility—and asked to comment on the confounders and the plausibility of five candidate instruments, considering both their relevance and threats to validity.
- Task 4: Reflection on LLMs as Co-Scientists. The expert reflected on the role of LLMs as collaborators in early-stage IV discovery, and whether such tools might augment, rather than replace, the theoretical reasoning of applied economists.

503 Appendix F. Detailed Definitions

504 Appendix F..1 Granger Causality

Granger causality is a statistical test used to determine whether one time series is useful in forecasting

another. Formally, for time-indexed data $\{v_{i,t}, v_{j,t}\}_{t=1}^T$, we test whether past values of v_i help predict

 v_i beyond what is possible using past values of v_i alone.

We define the null and alternative hypotheses as follows:

$$H_0: v_i \text{ does not Granger-cause } v_i$$
 (5)

$$H_1: v_i$$
 Granger-causes v_i (6)

This is operationalized by estimating and comparing the residual variances from two autoregressive models:

Restricted model (without v_i):

$$v_{j,t} = \alpha_0 + \sum_{k=1}^{p} \alpha_k v_{j,t-k} + \epsilon_t^{(r)}$$
 (7)

Unrestricted model (including lags of v_i):

$$v_{j,t} = \beta_0 + \sum_{k=1}^p \alpha_k v_{j,t-k} + \sum_{k=1}^p \gamma_k v_{i,t-k} + \epsilon_t^{(u)}$$
(8)

The null hypothesis corresponds to testing:

$$\gamma_1 = \gamma_2 = \dots = \gamma_p = 0 \tag{9}$$

If the unrestricted model significantly reduces the prediction error compared to the restricted model, we reject H_0 and conclude that v_i Granger-causes v_i .

515 Assumptions:

516

517

518

526

527

529

509

- Both time series are weakly stationary.
- The lag length p is appropriately selected.
- The model is correctly specified (linearity, no omitted variables).

519 Appendix F..2 ATE Estimation of IVs

- 520 When estimating causal effects using instrumental variables (IVs), we typically recover the Local
- 521 Average Treatment Effect (LATE), not the overall average treatment effect (ATE). This is because IV
- methods rely on *compliers*—units whose treatment status is affected by the instrument. As a result,
- 523 the estimated effect pertains only to this subpopulation.
- Formally, suppose we have an instrument Z, a treatment T, and an outcome Y. Under the potential outcomes framework, each unit i has:
 - $T_i(1)$ and $T_i(0)$: potential treatment values if $Z_i = 1$ or $Z_i = 0$
 - $Y_i(1)$ and $Y_i(0)$: potential outcomes under treatment or no treatment
- We define the following groups:
 - Compliers: $T_i(1) = 1, T_i(0) = 0$
- Never-takers: $T_i(1) = 0, T_i(0) = 0$
- Always-takers: $T_i(1) = 1, T_i(0) = 1$
- **Defiers:** $T_i(1) = 0$, $T_i(0) = 1$ (typically ruled out by the monotonicity assumption)

533 Key Assumptions for LATE:

534

537

538

- 1. **Relevance:** $\mathbb{E}[T|Z=1] \neq \mathbb{E}[T|Z=0]$ (instrument affects treatment)
- 2. **Independence:** $Z \perp \!\!\! \perp (Y(0), Y(1), T(0), T(1))$ (instrument is as good as randomly assigned)
 - 3. Exclusion Restriction: Z affects Y only through T (no direct effect on outcome)
 - 4. Monotonicity: $T_i(1) > T_i(0)$ for all i (no defiers)
- Under these assumptions, the LATE is identified as:

$$LATE = \frac{\mathbb{E}[Y|Z=1] - \mathbb{E}[Y|Z=0]}{\mathbb{E}[T|Z=1] - \mathbb{E}[T|Z=0]}$$

$$(10)$$

- This ratio represents the average causal effect of T on Y for compliers only.
- Two-Stage Least Squares (2SLS): To estimate LATE in practice, we use a two-stage regression procedure.
- In this appendix, we provide additional theoretical insights into the consistency metric introduced in the main text, highlighting its connection to the variance and bias properties of instrumental variable estimators.

546 Appendix F..3 Consistency as a Measure of Instrument Validity

Let Z be a candidate instrument used to estimate the causal effect β via

$$\hat{\beta}_{\text{IV}}^{(Z)} = \frac{\text{Cov}(Z, Y)}{\text{Cov}(Z, X)}.$$
(11)

Assuming Z satisfies the classical instrument validity conditions (relevance and exclusion), the IV estimator is consistent and unbiased in large samples. When multiple valid instruments Z_1, Z_2, \ldots, Z_m are available, their estimates $\hat{\beta}_{\text{IV}}^{(Z_i)}$ should converge to the true causal effect β as sample size grows, resulting in low pairwise differences:

$$\lim_{T \to \infty} \mathbb{E} \left| \hat{\beta}_{\text{IV}}^{(Z_i)} - \hat{\beta}_{\text{IV}}^{(Z_j)} \right| = 0, \quad \forall i, j.$$
 (12)

- Large observed discrepancies suggest violations of instrument validity, such as weak instrument bias or direct pathways from Z_i to Y that bypass X.
- Relation to Instrument Strength and Bias The variance of each IV estimate depends inversely on the strength of the instrument, quantified by $Cov(Z, X)^2$. Weak instruments induce greater variability, leading to increased disagreement between estimates from different instruments.
- Additionally, bias from invalid instruments inflates the expected pairwise difference. Formally, for instruments Z_i and Z_j , the expected squared difference decomposes as

$$\mathbb{E}\left[(\hat{\beta}_{\text{IV}}^{(Z_i)} - \hat{\beta}_{\text{IV}}^{(Z_j)})^2\right] = \underbrace{\text{Var}(\hat{\beta}_{\text{IV}}^{(Z_i)}) + \text{Var}(\hat{\beta}_{\text{IV}}^{(Z_j)})}_{\text{variance component}} + \underbrace{\left(\text{Bias}(Z_i) - \text{Bias}(Z_j)\right)^2}_{\text{bias component}}.$$
 (13)

- This decomposition illustrates how the consistency metric reflects both random variation and systematic bias in the set of instruments.
- Implications for Instrument Selection The normalized consistency score introduced in the main text effectively summarizes these properties by comparing observed discrepancies to a baseline derived from random (invalid) instruments. A low score implies both low variance and low bias among the instruments, supporting their joint validity.
- In practice, this metric can guide the selection and refinement of instruments by:

- Identifying instruments that cause high disagreement, which may be candidates for exclusion.
- Providing a quantitative measure to compare different instrument sets.
 - Complementing formal tests of instrument validity such as overidentification tests.

70 Appendix G. Rejected IVs

566

567

568

569

575

Treatment	Outcome	Rejected IVs
GDP	Conflict	Rainfall
BMI	SBP	MR
Church attendance	Crime	Rainy days
Turnout	Vote share	Rainfall
Protests	Prices	Rainfall

Table 4: Treatment-Outcome Pairs with Rejected Instruments

Appendix H. Gapminder preprocessing by IV Co-Scientist

Table 5 presents key preprocessing statistics for each variable pair from Gapminder, including the observed correlation between treatment and outcome variables and the corresponding sample sizes used in the analysis. These metrics provide context on the data quality and strength of associations before causal inference.

Treatment	Outcome	Correlation	Number of Data Points
GDP	Health	0.902	2784
Income	Carbon emissions	0.832	1790
Sanitation	Child mortality rate	-0.812	2578
Poverty	Cholesterol	-0.842	3568
Female literacy rate	Number of kids per female	-0.812	2504

Table 5: Preprocessing Summary: Correlation and Sample Size by Treatment–Outcome Pair

576 Appendix I. IVs generated by IV Co-Scientist

Table 6 summarizes the sets of accepted and rejected instrumental variables (IVs) for each treatment-outcome pair, as suggested by GPT-40. The accepted IVs represent those variables the model
deemed more plausible instruments after a critique stage, while the rejected IVs are those filtered out
due to likely violations of IV assumptions.

Treatment	Outcome	Accepted IVs	Rejected IVs
GDP	Health	 Distance to the port Global commodity prices 	Colonial legal-origin dummies Fertile land Historical settler-mortality rates
Income	Carbon emissions	 Industrial or resource endowments Policy reforms Trade in country 	 Distance to the equator Railroad network density
Sanitation	Child mortality rate	 Groundwater depth Sewerage investment 	 Sanitation subsidy rollout schedule Distance to health center Terrain
Poverty	Cholesterol	 Cash-transfer age cutoff State minimum wage 	 Childcare-program timing State EITC rate
Female literacy rate	Number of kids per female	 Number of female teachers Raised compulsory school-leaving age Introduction years of a girls-only scholarship program 	 Distance to school Historical density of missionary girls' schools (pre-independence) UI replacement rate

Table 6: Accepted and Rejected Instruments by Treatment–Outcome Pair

581 Appendix J. Prompts

HypothesisGenerator (Instrumental Variable)

You are an economist helping to identify causal relationships. Given the treatment variable $\{T\}$ and the outcome variable $\{Y\}$, please provide a list of 5 possible instrumental variables that could help estimate the causal effect of $\{T\}$ on $\{Y\}$. The context of this treatment-outcome pair is $\{Context\}$. These should be variables that influence $\{T\}$ but do not directly affect $\{Y\}$ except through $\{T\}$. Think step by step. Return your answer with Answer = [list of 5 IVs]

HypothesisGenerator (Confounder)

You are an economist helping to identify causal relationships. Given the treatment variable $\{T\}$ and the outcome variable $\{Y\}$, please provide a list of 5 possible confounding variables that might affect both $\{T\}$ and $\{Y\}$, potentially biasing the causal effect estimate. The context of this treatment-outcome pair is $\{Context\}$. Think step by step. Return your answer with Answer = [list of 5 confounders]

HypothesisGenerator (Independence)

You are an economist evaluating the validity of instrumental variables. Given the treatment variable $\{T\}$, outcome variable $\{Y\}$, a candidate instrumental variable $\{Z\}$, and a list of confounders $\{U_1,U_2,\ldots\}$, please assess the independence criteria i.e. $\{Z\}$ must be independent of any confounders that affect both $\{T\}$ and $\{Y\}$. Based on these definitions and the $\{Context\}$, please evaluate whether $\{Z\}$ is a valid instrument. Think step by step. Return your answer with Answer = [Valid]

584

583

HypothesisGenerator (Exclusion)

You are an economist evaluating the validity of instrumental variables. Given the treatment variable $\{T\}$, outcome variable $\{Y\}$, a candidate instrumental variable $\{Z\}$, please assess the exclusion criteria i.e. $\{Z\}$ affects the outcome $\{Y\}$ only through the treatment $\{T\}$, with no direct effect on $\{Y\}$. Based on these definitions and the $\{C\text{ontext}\}$, please evaluate whether $\{Z\}$ is a valid instrument. Think step by step. Return your answer with Answer = [Valid / Invalid].

585

ProxyHuman

You are a policy-minded economist tasked with identifying socio-economically meaningful causal questions. Given a candidate pair $\{(T, Y)\}$, assess whether: 1. The relationship is important or interesting i.e., is this a question researchers or policymakers would care about? 2. The pair is interpretable and policy-relevant in real-world socio-economic contexts. 3. The question could plausibly be studied using observational data. Avoid pairs that are too similar in meaning (e.g., literacy at ages $5\{10$ and literacy at ages $10\{15\}$). Think step by step, using the reasoning a social scientist or economist might apply when deciding whether to pursue this question.

586

CausalOracle

You are an economist reasoning about causal direction between two socio-economic variables. Given a variable pair (A, B) with a strong observed correlation, your task is to determine the likely causal relationship. Please evaluate: 1. Is it more plausible that A causes B? 2. Is it more plausible that B causes A? 3. Could the relationship be bidirectional? 4. Or is the correlation likely driven by confounding or coincidence, with no direct causal link? Use real-world knowledge and reasoning as an economist to assess plausibility. Think step by step. Return your answer as: Answer = [1 / 2 / 3 / 4].

587