
Read the Paper, Write the Code: Agentic Reproduction of Social-Science Results

Benjamin Kohler
ETH Zurich

David Zollikofer
ETH Zurich

Johanna Einsiedler
University of Basel

Alexander Hoyle*
ETH Zurich

Elliott Ash*
ETH Zurich

Abstract

Recent work has used LLM agents to reproduce empirical social science results with access to both the data and code. We broaden this scope by asking: Can they reproduce results given only a paper’s methods description and original data? We develop an agentic reproduction system that extracts structured methods descriptions from papers, runs reimplementations under strict information isolation—agents never see the original code, results, or paper—and enables deterministic, cell-level comparison of reproduced outputs to the original results. An error attribution step traces discrepancies through the system chain to identify root causes. Evaluating four agent scaffolds and four LLMs on 48 papers with human-verified reproducibility, we find that agents can largely recover published results, but performance varies substantially between models, scaffolds, and papers. Root cause analysis reveals that failures stem both from agent errors and from underspecification in the papers themselves.

1 Introduction

Agentic systems are increasingly able to autonomously generate, debug, and execute end-to-end research pipelines across a broad range of domains (Lu et al., 2024; Si et al., 2025; Gandhi et al., 2025; Schmidgall et al., 2025; Sun et al., 2025). These capabilities have led to the application of agents to the task of *scientific reproduction*, as well as to novel considerations of replication-oriented publication pipelines (Brodeur & Barbarioli, 2025; Fishman & Sekeres, 2025). In the empirical social sciences specifically, recent work uses agents to execute the preexisting code associated with papers and to assess their level of reproducibility (Hu et al., 2025; Shah et al., 2026; Xu & Yang, 2026).

Yet publications, not code, communicate scientific research and serve as the source of truth. An open question, then, is whether reproduction of research results is contingent on direct access to the analysis code. Or whether, as presumably intended, the paper’s methods description is sufficiently detailed to allow reproduction through re-implementation of the code from scratch.

This paper makes progress on this question by building an agentic system to reproduce research results given a paper’s method description and associated data, but without the analysis code. Our pipeline extracts structured representations of the research methods and results, where specific numbers are masked. The agent is then tasked with reproducing the masked values by implementing the analysis pipeline from scratch, allowing for deterministic benchmarking.

*Equal supervision. Author contact: Kohler, benjamin.kohler@gess.ethz.ch; Zollikofer, david.zollikofer@gess.ethz.ch; Einsiedler, johanna.einsiedler@unibas.ch; Hoyle, alexander.hoyle@ai.ethz.ch; Ash, ashe@ethz.ch. Code at github.com/benjamin-kohler/social_science_replicability.

We apply the pipeline to a dataset of 48 social science papers that have been human-verified as reproducible (Institute for Replication, 2024). Overall, frontier LLM agents are quite successful at re-implementing analysis code of social science papers. For the three best-performing agents, reproduced and original coefficients agree on the sign over 85% of the time, and they are within the 95% confidence intervals over 70% of the time.

Performance varies a lot across LLMs and across agent scaffolds. The best performing agent in our comparisons is GPT-5.4 using the OpenCode scaffold, out-performing Claude Code (Opus 4.6), GPT-5.3 Codex, and GLM-5 using the OpenCode scaffold. The scaffold-model interaction matters a lot, with GPT-5.4 on OpenCode outperforming GPT-5.4 on the Codex CLI or on mini-SWE-Agent. In supporting analysis, we show that OpenCode GPT-5.4’s better performance is due to much greater token usage, coming at the cost of more expensive API calls and longer analysis runs.

Our reproduction pipeline enables a detailed diagnosis of failure modes to help identify the source of incorrect estimates. We find that errors mostly come from papers’ under-specification of methods and—to a lesser extent—from agents’ misunderstanding of or noncompliance with the methods description.

These methods and findings expand on the recent literature on automated scientific reproducibility. While important gaps remain to be closed, LLM agents can increasingly *read the paper and write the code*. These systems could play an important role in addressing the well-recognized reproduction gaps in empirical social science (Ioannidis, 2005; Nosek et al., 2012; Brodeur et al., 2024).

2 Background on agentic reproducibility

Independent verification is central to the scientific method (Nosek et al., 2022). Recent crowd-sourced reproduction efforts have revealed substantial variation in reproducibility across social science disciplines. Yet such work requires considerable resources—hundreds of trained researchers to collectively assess only a negligible fraction of the published literature (Brodeur et al., 2026; Miske et al., 2026; Tyner et al., 2026).

A helpful framework from Dreber & Johannesson (2025) organizes different forms and components of reproducibility and replicability. *Reproducibility* tests whether the results and conclusions of original studies can be reproduced based on the *same data* used in the original studies. *Replicability* concerns the validity of results when similar methods are applied to *new* data. This paper concerns reproducibility.

Reproducibility can be further decomposed into three types: (1) re-running the original authors’ code on identical data to regenerate results, (2) reproducing results using only the information reported in the paper and the same data, and (3) testing whether results are robust to reasonable alternative analytical decisions applied to the same data. Our work falls under (2), re-implementation.

This paper contributes to ongoing research exploring the deployment of LLM agents for reproduction and replication, summarized in Table 1.¹ Existing benchmarks in the ML domain evaluate agents’ ability to reproduce code architectures and results. In social science, prior works—Hu et al. (2025), Zhang et al. (2026), Shah et al. (2026), and Xu & Yang (2026)—benchmark reproducibility while granting agents access to the original code. Our contribution differs in that we reproduce social science results from the methods and data alone, without access to any author-provided code. In Nguyen et al. (2026), the agents try

¹We exclude benchmarks of code implementation quality rather than reproducibility (Huang et al., 2023; Tian et al., 2024; Hua et al., 2025).

System	Domain	Has Code	New Data	Evaluation
CORE-Bench (Siegel et al., 2024)	Multi	✓		Det.
PaperBench (Starace et al., 2025)	ML			LLM
Paper2Code (Seo et al., 2026)	ML			LLM
AutoReproduce (Zhao et al., 2025)	ML			Both
FIRE-Bench (Wang et al., 2026)	ML			LLM
REPRO-Bench (Hu et al., 2025)	Soc. sci.	✓		LLM
PaperRepro (Zhang et al., 2026)	Soc. sci.	✓		LLM
Shah et al. (2026)	Soc. sci.	✓		Det.
Xu & Yang (2026)	Soc. sci.	✓		Det.
ReplicatorBench (Nguyen et al., 2026)	Soc. sci.	✓	✓	LLM
This work	Soc. sci.			Det.

Table 1. Automated systems for scientific reproducibility. Columns indicate disciplinary domain (multidisciplinary, machine learning, or social science), whether the agent has access to the code, whether the agent replicates on new data, and the evaluation approach (deterministic, LLM-as-judge, or both).

to replicate (rather than reproduce) a paper by retrieving new data and testing the same hypothesis.²

Unlike most prior work, which relies on LLM-as-judge evaluation, our evaluation is deterministic, comparing reproduced outputs directly to original values with adjustments for statistical significance based on ground-truth standard errors. This approach may be more reliable, since LLMs are not well-validated for evaluation tasks (Li et al., 2024).

3 A pipeline for paper-derived reproduction

We develop a multi-step pipeline to evaluate whether LLM agents can reproduce social science results. The pipeline tasks agents with re-implementing code from methods descriptions extracted from the paper, using the associated data (see Figure 1). Separating the pipeline into distinct steps minimizes information leakage of both existing numerical results and original code to the agents. The four steps are: (1) extracting data, methods, and results from the paper and associated files; (2) re-implementing the code based on the data, methods, and results table templates; (3) evaluating estimated results against the original results; and (4) diagnosing and explaining deviations.

3.1 Step 1: Extraction

The first step extracts structured information from the paper and its reproduction package in three sub-processes (for the exact prompt see Appendix C.1).

Method extraction. We provide an LLM with the complete PDF of the paper and prompt it to extract a complete methods description, similar to the approach taken by Seo et al. (2026) for machine learning papers. The methods description includes the main research question, the study setting (e.g., country and time period), data descriptions (e.g., data sources and a broad summary of each dataset), general data manipulations (e.g., filtering

²In one of Nguyen et al.’s (2026) specifications, the agent does not have access to the analysis code from the replicated paper, but it still has access to the complete original PDF and parsed results—the answers it aims to replicate. We remove such information in our approach because it can cause leakage—in preliminary versions of our setup where models had access to results, agents would often copy them directly.

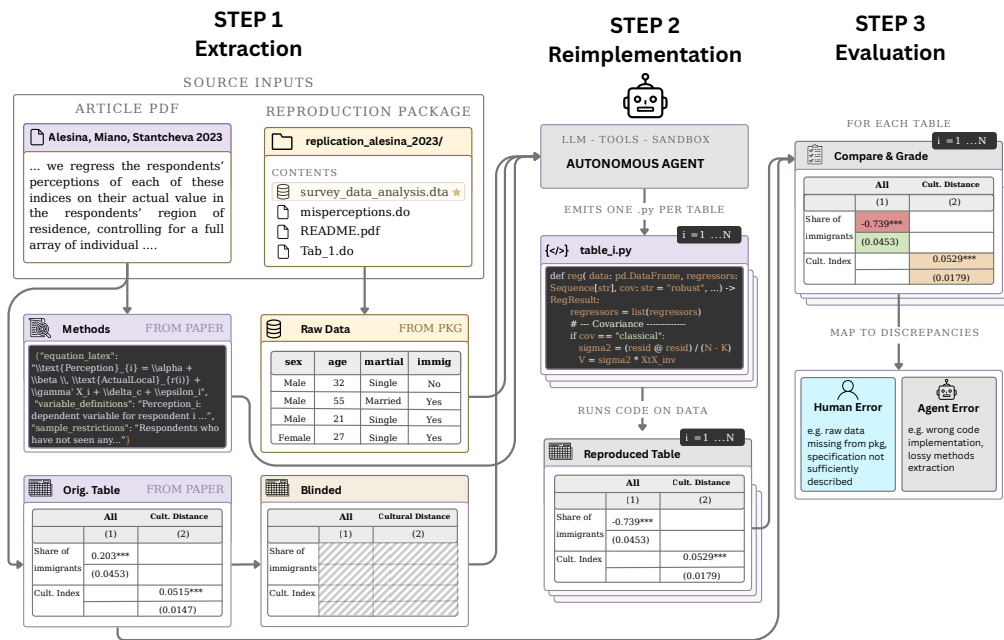


Figure 1. Overview of the pipeline for replicating empirical results. Step 1 (Extraction): Paper and replication package are parsed to extract the methods, original data, and a blinded version of the original results table. Step 2 (Reimplementation): An autonomous LLM agent, equipped with tools and a sandboxed execution environment, emits one Python script per table and runs it on the original data to produce a reproduced table. Step 3 (Evaluation): Each reproduced table is compared cell-by-cell against the original, and divergences are traced back to human errors (e.g., missing data, underspecified methods in the original package) and agent errors (e.g., incorrect code, lossy extraction of the methods). Code at github.com/benjamin-kohler/social_science_replicability.

steps or variable construction), and detailed descriptions of each results table without the numerical contents (e.g., name, caption, row and column labels, table-specific data filters, and model specifications such as regression formulas). The model is explicitly instructed not to include any numerical results or conclusions, verified by checking for any numerals in the table description.

This step produces a structured representation of the paper’s method without leaking results or code to the reimplementing agent (Section 3.2). Compared to providing the original paper directly (as is done in [Nguyen et al. 2026](#)), this step helps prevent the agent from hard-coding results or re-iterating code until exact results are reached.

The quality of the methods extraction influences downstream results. An incomplete or erroneous extraction will result in a failed replication. In the error source diagnosis below, we show that incomplete extraction is occasionally responsible for failed reproductions. While we found that overall, the extracted methods descriptions were of high quality, there is likely scope for further improvements in the quality of method extraction.

Results extraction and blinding. The next step is to extract the original results from tables.³ A dedicated LLM call retrieves the numerical values from all tables in the paper.⁴ The model returns a structured representation of the table in which cells are identified by both row and column indices as well as their label names. For each cell, the model records the content as a string. It then decomposes it into several components: a numerical value, a classification of the metric type (e.g., coefficient or number of observations), the number of significance stars, and—for statistics such as standard errors or p-values—a reference to the corresponding coefficient.

To construct a template for the reproduction agents to fill in, we “blind” a copy of this structured table by setting all cell values and significance stars to null.

Table parsing uses GPT-5-mini applied to the PDF rendered as images, which we found to outperform extraction from machine-readable PDF text. We manually validate table extraction on a random sample of five papers (24 tables in total). The model extracted information perfectly for all portrait-oriented tables and produced occasional minor numerical errors in the second or third decimal place for landscape-oriented or multi-page tables.

Data extraction. We instruct an LLM (GPT-5-mini) to identify the data files required for replication from the full directory tree and README of each provided reproduction package. The model is instructed to limit the datasets to the least pre-processed version available, i.e., the minimum viable set to perform reproduction. The model also classifies whether the package contains sufficient data to fully reproduce the paper, or whether data were partially or entirely missing—for example, because they are proprietary, subject to privacy restrictions, or require independent collection. Papers with no usable data are skipped.

3.2 Step 2: Reimplementation

The extracted method, table templates, and original data serve as the sole inputs for the reproduction agents, following [Starace et al. \(2025\)](#) and [Seo et al. \(2026\)](#). Agents are instructed to read the extracted methods and table templates, explore the data, and write a separate

³Empirical social science results can take the form of either tables or figures. Our pipeline handles both tables and figures, but we restrict our core analysis to tables because the results can be compared deterministically. Figure-level comparisons require vision-language models whose reliability for this task remains uncertain ([Wang et al., 2024](#); [Tang et al., 2025](#); [Hou et al., 2026](#), *inter alia*). Supporting figure-level results are reported in Appendix A.2.

⁴We include all tables from the paper’s published journal PDF. Tables in separate online appendices are not included. In the results below, about 8% of tables (from 18 papers) are not included in the analysis due to extraction issues.

Python script for each table that fills in the corresponding template (see Appendix C.2 for the prompt used).

The agents are restricted from accessing any file paths outside their workspace and are prohibited from retrieving the original reproduction package or paper via web tools. Each agent operates in an isolated environment containing only the task description, the extracted methods, the output templates, and a symlink to the data. The original reproduction package and paper PDF remain on the same machine but outside this boundary. These restrictions are necessary because agents could otherwise simplify the task by directly accessing the original code or results.

We implement a two-stage audit pipeline to verify compliance with the guardrail (Appendix B). The guardrail audit includes a deterministic regex scan, which classifies all accessed file paths and URLs as allowed or forbidden, as well as a GPT-5.4-mini review that rates each run from clean to severe violation. Suspicious runs are manually inspected and rerun if leakage is detected. In parallel, a hardcoding audit checks whether agents output statistical results as numeric literals without a computation path from the data. These two audits provide complementary signals: the guardrail audit identifies how forbidden information may have entered the pipeline, while the hardcoding audit detects outputs not derived from genuine computation. Together, they indicate that the reported results reflect genuine reimplementations rather than retrieval or memorization of pre-existing outputs.

3.3 Step 3: Evaluation

Since both the agent-produced results and the original results share the same structured template, we can directly compare outputs at the cell level. For all numerical values, we check the sign agreement and compute the percentage-point difference. For regression coefficients, we additionally compute the difference scaled by the ground-truth standard error.

Each numerical value receives a letter grade based on the percent absolute deviation from the original value. If the sign matches, we assign A (0%-2%), B (2%-20%), C (20%-40%), D (40%-60%), or E (60+%). If the sign doesn't match, we assign E. If the estimate is not produced or is not a number, we assign F (see Appendix A.3 for a few more details on grades).

The letter grades have two advantages over the statistical scores. First, assigning E for 60+% discrepancies is more robust to outlier estimates. Second, multiple criteria like the sign match can be combined into one grading scheme.

Results are aggregated first by computing the average grade across all cells within a table, and then by computing the average table grade within a paper. Binning into grades before aggregation improves interpretability and reduces the influence of outliers. Unlike some of the prior work, our approach avoids the ambiguity and opacity of LLM-judges. In this regard, it is closest to Xu & Yang (2026), who also measure exact reproduction of coefficients.

3.4 Step 4: Explanation

The explanation step produces interpretable diagnostics of discrepancies between the original and reproduced results. An LLM agent (GPT-5.4 with Codex CLI) is prompted to perform the analysis as follows. The input is the set of observed cell-level numerical failures (all cells with a score of B or below). For each failing output, the agent is tasked to locate the relevant code in both the original reproduction package and the agent-generated Python code. It should then identify and describe the discrepancy causing the failure.

For each reported discrepancy, an LLM auditor runs consistency checks comparing the various input and output files. The auditor can identify four main types of errors. In the bucket of *Human Error*, we have *Missing Data* and *Paper vs. Code* (the original paper did not properly or completely describe the code or contradicts it). Under *Agent Error*, we have *Paper vs. Methods Extraction* (the methods extraction pipeline did not properly or completely describe the methods) or *Method Extraction vs. Agent* (the agent did not properly follow the extracted methods description).

4 Application to I4Replication data

4.1 Data

Our analysis is based on papers with verified reproducibility, enabling meaningful comparisons across reproduction results. We draw on the work of *I4Replication*, an institute coordinating a large-scale research community effort to reproduce social science papers ([Institute for Replication, 2024](#)). Their process combines systematic review of reproduction packages, step-by-step code execution, and direct exchange with the original authors, resulting in a detailed reproducibility report and classification for each paper. We collect the PDFs and reproduction packages of all papers that I4Replication has classified as fully reproducible, leading to a total dataset of 48 papers (see Appendix Table A1).

Appendix Table A3 lists all the papers with associated journal, title, main programming language, and number of lines of code. Table 2 Panel A tabulates the counts by journal. The papers are from a mix of well-regarded economics and political science journals. Summary stats on the code are shown in Appendix Table A2. The original analysis code is mostly programmed in Stata (54%) or R (27.1%), with zero papers in Python. The packages have an average of 5,324 lines of code.

Summary statistics on the extracted paper elements are reported in Table 2 Panel B. The system extracted 222 ground-truth tables containing 14,214 table cells. Of these, about one-third (5,149) contain regression coefficients. About four-fifths (4,253) of coefficients are accompanied by a standard error, while one-sixth (779) have p-values. Other frequently observed statistics are the number of observations and R-squared.

A potential issue with this application is data leakage: tested LLMs may have been exposed to the manuscripts or code during pre-training. To probe the relevance of leakage, we provide an auxiliary analysis reproducing results from papers published after the model knowledge cutoff, which is August 2025 for both GPT 5.3-Codex/5.4 ([OpenAI, 2026](#)) and Claude Opus 4.6 ([Anthropic, 2026](#)). Specifically, we sample ten papers (with reproduction packages) from the *Economic Journal* (EJ), with five papers published *before* and five papers published *after* August 2025. While these papers have not yet been verified as reproducible by I4Replication, they have been verified by EJ’s data editor team, which carefully vets reproduction packages and verifies that they produce the same outputs as the paper.

4.2 Agent systems

The LLM agents tested are GPT-5.4, GPT-5.3 Codex, Claude Opus 4.6, and the open-weights GLM-5 ([Team et al., 2026](#)). We evaluate two proprietary agent scaffolds—Claude Code and Codex—and two open-source alternatives—mini-SWE-agent ([Yang et al., 2024](#)) and OpenCode. For the open-source scaffolds, we compare a proprietary model (GPT-5.4) and an open-weights alternative (GLM-5). Figure A.8 illustrates some sample code segments generated by the agents.

Before evaluating success in reproduction, a preliminary question is whether these agent systems can even generate any results at all. Table 2 Panel C summarizes the completion

Panel A: Papers by journal and discipline (final sample, $N = 48$)

Journal	Discipline	N papers
<i>American Journal of Political Science (AJPS)</i>	Political Science	11
<i>Economic Journal (EJ)</i>	Economics	9
<i>American Economic Review (AER)</i>	Economics	8
<i>American Economic Journal: Economic Policy (AEJ:Pol)</i>	Economics	5
<i>Journal of Politics (JOP)</i>	Political Science	5
<i>American Economic Journal: Applied Economics (AEJ:AE)</i>	Economics	3
<i>American Political Science Review (APSR)</i>	Political Science	3
<i>Review of Economic Studies (REStud)</i>	Economics	2
<i>American Economic Journal: Macroeconomics (AEJ:Macro)</i>	Economics	1
<i>Quarterly Journal of Economics (QJE)</i>	Economics	1
Total		48

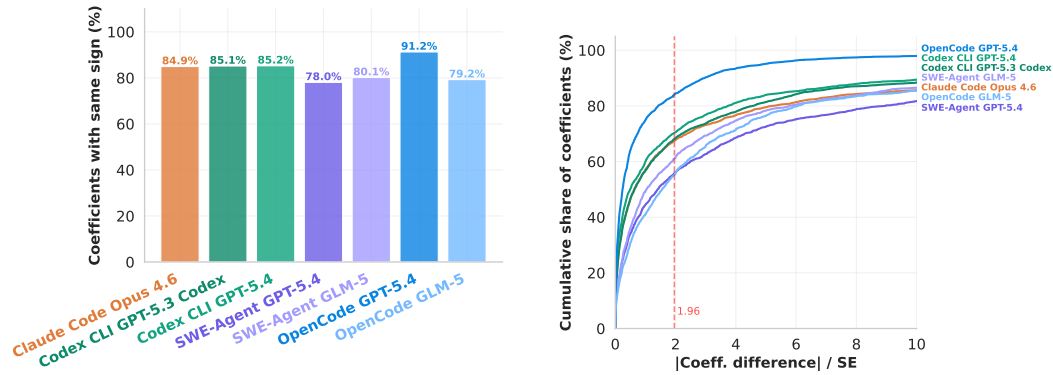
Panel B: Extracted Elements in Original Papers (per paper, $N = 48$)

	Total	Mean	SD	Min	Max
Tables	222	4.6	2.8	1	11
Cells (present)	14,214	296.1	261.4	20	1433
Coefficients	5,149	107.3	93.5	8	384
Standard errors	4,253	88.6	85.3	0	319
p-values	779	16.2	32.5	0	127
t-statistics	10	0.2	1.4	0	10
Confidence intervals	112	2.3	16.2	0	112
R-squared	590	12.3	17.9	0	58
N observations	1,701	35.4	86.2	0	602
F-statistics	121	2.5	6.2	0	30
Other numeric	1,607	33.5	53.6	0	253

Panel C: Completion Counts and Rates (%) by Paper Element and Agent System

Scaffold LLM	Claude Code		Codex CLI		SWE-Agent		OpenCode		Average
	Opus 4.6	GPT-5.3	Codex	GPT-5.4	GPT-5.4	GLM-5	GPT-5.4	GLM-5	
Papers	48 (100%)	45 (94%)	45 (94%)	46 (96%)	45 (94%)	45 (94%)	45 (94%)	44 (92%)	95%
Tables	213 (96%)	215 (97%)	214 (96%)	204 (92%)	182 (82%)	207 (93%)	213 (96%)		95%
Cells (Overall)	10,276 (72%)	9,566 (67%)	9,615 (68%)	8,101 (57%)	7,374 (52%)	9,096 (64%)	7,648 (54%)		62%
Coefficients	4,355 (85%)	4,315 (84%)	4,395 (85%)	4,113 (80%)	3,907 (76%)	4,180 (81%)	4,326 (84%)		82%
Standard errors	3,414 (80%)	3,623 (85%)	3,541 (83%)	3,257 (77%)	2,955 (69%)	3,458 (81%)	3,443 (81%)		80%
p-values	692 (89%)	590 (76%)	697 (89%)	617 (79%)	663 (85%)	501 (64%)	655 (84%)		81%
t-statistics	10 (100%)	10 (100%)	10 (100%)	10 (100%)	10 (100%)	10 (100%)	10 (100%)		100%
Confidence intervals	58 (52%)	54 (48%)	54 (48%)	54 (48%)	54 (48%)	54 (48%)	54 (48%)		49%
R-squared	498 (84%)	524 (89%)	488 (83%)	506 (86%)	424 (72%)	428 (73%)	502 (85%)		82%
N observations	1,267 (74%)	1,402 (82%)	1,397 (82%)	1,315 (77%)	858 (50%)	1,341 (79%)	1,383 (81%)		75%
F-statistics	108 (89%)	109 (90%)	105 (87%)	111 (92%)	62 (51%)	109 (90%)	111 (92%)		84%
Other numeric	1,175 (73%)	1,326 (83%)	1,248 (78%)	1,261 (78%)	1,036 (64%)	1,296 (81%)	1,240 (77%)		76%

Table 2. Descriptive overview Panel A: number of included papers by journal. Panel B: per-paper summary statistics on the original outputs. Panel C: reproductions, cell counts given as absolute totals with the percentage of originals in parentheses; Average is the mean of the shares across the 7 approach–model combinations.



(a) Share of reproduced coefficients with the same sign as the original paper.

(b) Cumulative distribution of $|\text{coefficient difference}| / \text{SE}$ across approaches. The standard error is the original standard error of the coefficient reported in the paper. The dashed red line marks 1.96 (95% CI threshold).

Figure 2. Performance metrics of agent-reproduced coefficients. *Coefficients are identified via an LLM table extraction pipeline that classifies all numerical values into one of multiple statistic types. Agentic systems can reliably recover a large share of the directional findings from social science publications given the publication text and data alone. All metrics are based on reproduced results excluding missing coefficients. Panel (a): The best performing agent recovers the correct sign over 90% of the time. Panel (b): The reproduced coefficients of the best performing agent are within the 95% confidence interval (based on the ground-truth standard error) over 80% of the time.*

rates according to various features of the paper results. The agents deliver usable results for at least 92% (and up to 100%) of papers, 82% (up to 97%) of tables, and 52% (up to 72%) of cells. The bottom part of the table shows completion across various types of cells, with the most important ones being coefficients (82% completion on average) and standard errors (80% completion on average). In terms of overall completion, the best agent system is Claude Opus 4.6 and the worst is SWE-Agent GLM-5.

5 Results: How well do AI agents reproduce papers without their code?

5.1 Main results

This section presents the main results on reproducing the I4Replication papers. We start at the cell (estimate) level and proceed to table- and paper-level results.

Cell-level. Figure 2 reports cell-level results for regression coefficients, which are the most prevalent statistic in tables. First, Figure 2a reports the share of reproduced estimates that match the sign of the original estimates. Across all models and scaffolds, reproduced coefficients match the sign for the large majority of observations, ranging from 78% (SWE-Agent GPT-5.4) to 91% (OpenCode GPT-5.4). In all cases, that makes a significant lift over the naive “guess positive”, which would give 68% accuracy. Note that the numbers in Figure 2a ignore “missing” coefficients that were not generated. For comparison, Appendix Figure A3 includes missing coefficients in the denominator.

Next, Figure 2b presents cumulative distributions (CDFs) of the absolute difference of the original and reproduced coefficients, divided by the ground-truth standard error. This normalized difference measures the discrepancy between original and reproduced coefficients in statistically meaningful units (same units as the Wald t-statistic). For these curves, better performance is indicated by hugging the top-left corner. The value of the CDF at $x = 1.96$

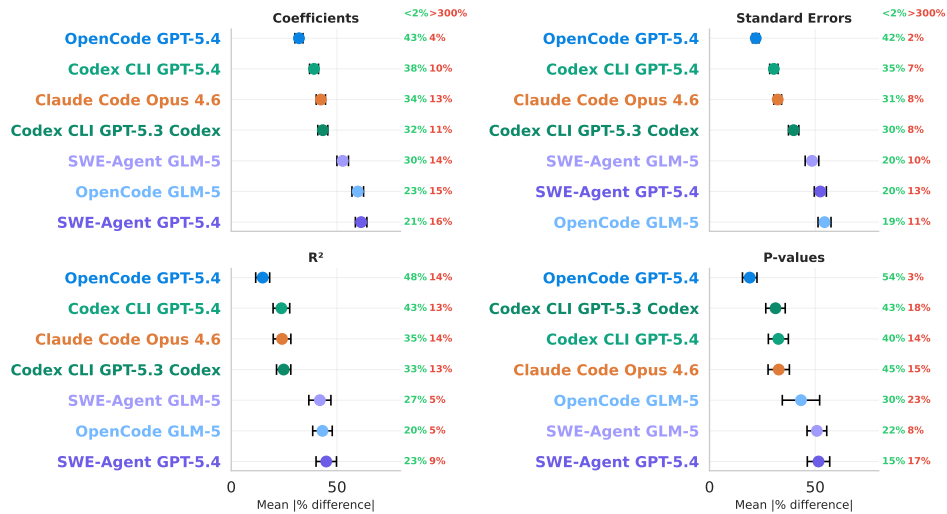


Figure 3. Mean percentage differences between original and reproduced, by statistic type. The strongest models can reproduce more than 40% of all coefficients and standard errors exactly. Note that these statistics are only produced when the original table contains them, and therefore the number of observations per table varies. We drop observations beyond a 300% difference, as these outliers frequently are due to scaling issues. Green and red numbers next to the plot indicate both the share of these large outliers (red) and the share that are (close to) perfectly reproduced (green).

indicates the share of reproduced estimates which lie within the 95% confidence interval (CI), meaning they are not statistically different from the original estimates. Here, we see that the reproduced coefficient falls within the 95% CI for more than half of all reproduced values, even in the worst-performing case. For the best model (OpenCode GPT 5.4), over 80% of reproduced values are within the 95% CI.

The CDFs further reveal that a non-negligible share of reproduced coefficients deviate substantially from the originals. An informal inspection of these cases suggests that while some are attributable to agent behavior (e.g., selecting the wrong dataset from among several similar candidates), the majority stem from unit scaling issues not accounted for by our comparison system (e.g., the original authors report values in dollars, whereas the dataset is denominated in cents).

Figure 3 measures reproduction error as the absolute percentage-point difference between the original and reproduced statistic, with additional heterogeneity reported by type of statistics (coefficient, standard error, R-squared, and p-value). Across all statistic types, the best agent is OpenCode GPT 5.4. The worst agent is SWE-Agent GPT-5.4 (3/4 stats) or OpenCode GLM-5 (1/4 stats). Codex GPT-5.4, Codex GPT-5.3, and Claude Code Opus 4.6 perform similarly, forming a second tier behind OpenCode GPT-5.4. SWE-Agent performs the worst among the scaffolds, even with GPT-5.4 as the underlying model.

Table- and paper-level. Figure 4 reports results with letter grades, aggregated by table (panel a) and paper (panel b). Converting numerical values to grades mitigates the influence of outliers when comparing and aggregating results. Overall, the results are consistent when looking at tables and papers. The table-level aggregation highlights performance differences across agents: the best-performing model-scaffold combination reproduces twice as many tables near-perfectly (grade A) compared to the worst-performing one.

Both the table- and paper-level aggregations indicate that achieving accurate and reliable reproduction beyond individual cells remains challenging for all agents. Agent rankings become less differentiated at the paper level, suggesting that some agents tend to perform

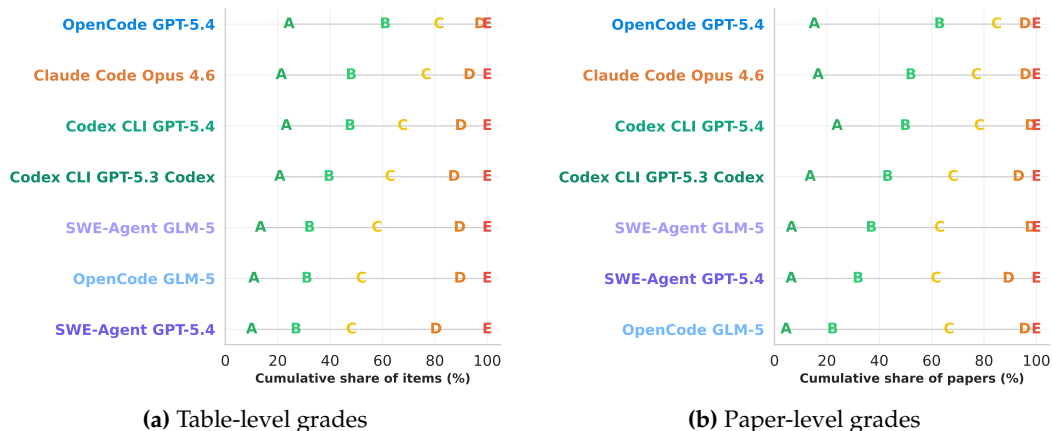


Figure 4. Grade comparison by aggregation level and agent. Grades are assigned to percentage differences to reproduced and original results following the rubric in Table A5 and then aggregated with equal weights. All reproduced numbers are rounded to the number of digits reported in the original papers. The best scaffolds and models reproduce more than a fifth of tables almost exactly (average cell within 2%, an A grade), and more than 60% within 20% (a grade B); Comparable, but slightly smaller proportions hold at the level of the entire paper. Observations are ranked by their share of A and B grades. Equivalent distributions including the F-grades can be found in Appendix Figure A2.

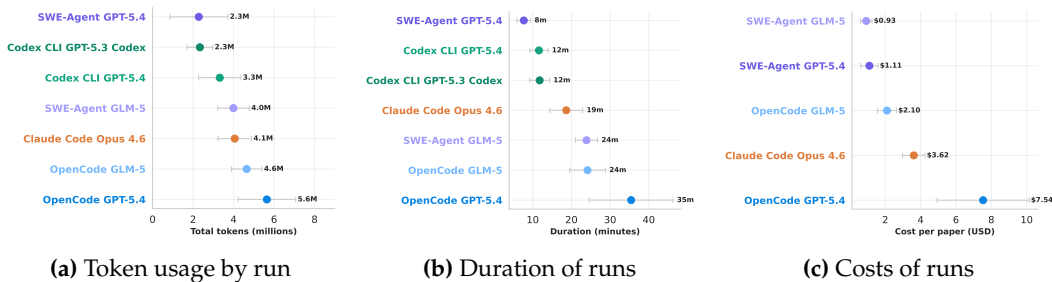


Figure 5. Token usage, run durations, and run costs, by agent. Error bars indicate 95% confidence intervals. Token usage includes cached and uncached input and output tokens. Codex CLI (through subscription) did not report usage costs. OpenCode in combination with GPT 5.4 consumes more tokens, takes longer, and costs more.

consistently well across all tables once they successfully reproduce one (e.g., Codex with GPT-5.4), whereas others solve individual tables more independently, leading to greater within-paper performance variance.

Appendix Figure A6 shows table results separately by the empirical function of the table – main results, mechanism analysis, robustness checks, or descriptive statistics. The results show that descriptive statistics are easier to reproduce than the results tables, which demonstrate similar performance across main, mechanisms, and robustness. Appendix A.6 provides a number of other supporting results on correlates of reproduction performance.

5.2 Variation in agentic effort

To better understand performance differences across agents, we examine variation in computational effort, measured by token usage, runtime, and cost. Figure 5 shows substantial heterogeneity along all three dimensions.

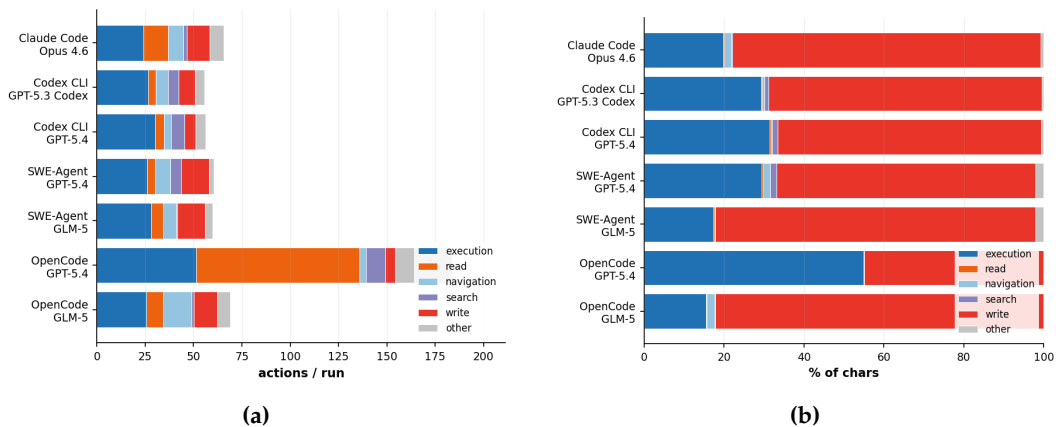


Figure 6. Effort allocation across tasks. The figures show the number of tool-call actions taken by the agents (left) and the relative volume of text (number of characters) emitted (right), colored by function category.

The best-performing agent, OpenCode GPT-5.4, consistently operates at a higher level of effort. It consumes many more tokens and takes much more time per run than competing approaches, and it is also the most costly configuration. This pattern suggests that its superior performance is at least partly driven by greater computational investment: more extensive exploration, longer reasoning chains, or more iterations over the data and code. By contrast, lower-performing agents tend to operate under tighter computational budgets, which may limit their ability to resolve ambiguities or recover from intermediate errors.

Token usage statistics give us an indication of the overall effort on a task but not how it is spent. To better understand the composition of this effort, we use an LLM to classify each agent action into one of six categories (execution, reading, navigation, search, writing, or other) and count both the number of actions and the volume of text (in characters) produced in each action (details in Appendix B.2).

Figure 6a shows the breakdown. For most models, the main tool call is execution, followed by a mix of reading and writing, and then other tools. However, OpenCode GPT-5.4 makes more than twice as many tool calls as other setups, spending more calls on execution and many more calls on reading files. When weighting the tool calls by volume of text (right panel), however, the picture looks very different. Most tokens are spent on writing, with almost all of the remaining tokens spent on execution. As a share of tokens, reading, searching, navigation, and other are negligible tasks. Qualitative analysis reveals that while GPT-5.4 mostly favors inline Python execution via heredocs, GLM-5 writes scripts to disk before execution, explaining the larger fraction of characters spent on writing files.

These results highlight a key trade-off in agentic reproduction. Performance gains are a function of the effort (tokens) expended during execution. In this sense, the strong performance of OpenCode GPT-5.4 reflects both effective scaffold-model alignment and a willingness to spend more tokens and time on each task. This finding is consistent with prior work emphasizing the importance of scaffold-model interaction for agentic performance (Cohn et al., 2026), and suggests that differences in observed accuracy may partly reflect differences in implicit compute budgets rather than purely differences in capability.

5.3 Exploration of error sources

Where do these errors originate? Figure 7 summarizes the root causes identified by our discrepancy pipeline across model-scaffold combinations, grouped into five categories:

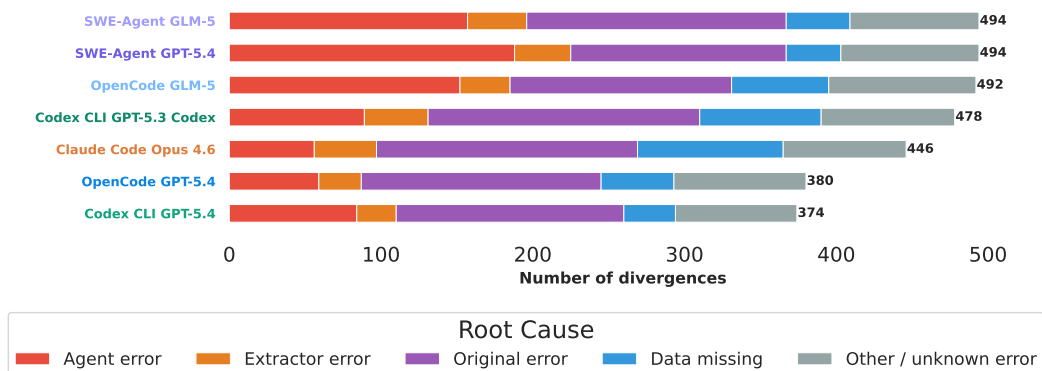


Figure 7. Error analysis pipeline: Number of divergences by error source. The error pipeline identifies divergences between the original and reproduced code for tables that score worse than overall grade B. These divergences are described and classified into categories. The agent then traces the root cause of the divergence. **Agent error:** the agent makes an error in mapping the methods to reproduced code. **Extractor error:** the initial extraction from paper to methods description is faulty. **Original error:** the paper under-specifies or mis-specifies the original code. **Data missing:** data is missing in the reproduction package. Appendix A.7 has more detail.

agent error, extractor error, original error, data missing, and other/unknown. As shown, the large majority of divergences—over three-fourths—can be traced to a specific, interpretable source.

A substantial share of discrepancies are outside the reproducing agent’s control. First, the largest share of divergences stems from *original errors*, indicating mismatches between the paper and the underlying code. In other words, the methods descriptions in the paper are often insufficiently precise to enable faithful reimplementations. Another substantial share of errors comes from missing data. Finally, a number of discrepancies arise from errors in the extraction pipeline – done in our shared data construction step, before the agent systems start their work.

Agent errors constitute the second largest category of discrepancies. For the strongest agents, the share of agent errors declines markedly, becoming a relatively minor component of total discrepancies. A larger relative fraction of divergences is attributed to human sources — either to underspecification in the original paper or to missing data. This pattern suggests that limitations in how methods are documented, rather than agent capability alone, are a primary barrier to automated reproducibility.

There is substantial variation across agents in the number of discrepancies attributed to human sources, such as paper underspecification. This variation does not imply that some agents encounter fewer human-caused issues, but rather that they handle these ambiguities differently. Even when the underlying discrepancy originates from incomplete or unclear reporting in the paper, agents may still reproduce the correct result if their implicit assumptions happen to align with the original implementation. As a result, the observed number of human-attributed discrepancies reflects both the presence of underspecification and the agent’s ability to resolve it correctly.

Table 3 illustrates this mechanism. In all three cases, the paper omits key details—such as how party affiliation is coded or which F-statistic is used—forcing the agent to infer the correct implementation. Different agents make different assumptions: for instance, one agent incorrectly maps party codes or relies on a default statistical routine, while another reconstructs the intended specification more accurately and recovers results closer to the original. These examples show that the same human-caused ambiguity can either lead

	10.1093/ej/ueab069 — Table 3a		10.1257/pol.20200559 — Table 2		10.1093/ej/ueab102 — Table 3	
	Behavior	Gr.	Behavior	Gr.	Behavior	Gr.
Claude code behavior	Agent assumes id for democrats is party=100	C	Agent assumes that codes A1 and A2 indicate democrat support	D	Uses <code>linearmodels.f.stat</code> diagnostics.	D
OpenCode 5.4 behavior	Agent assumes id for democrats is party=200	B	Agent assumes that codes A4 and A5 indicate democrat support	B	Manually reconstructs the (Kleibergen-Paap) Wald F-statistic	B
Correct behavior	Democrats id is party=200		A4 and A5 indicates democrat support		Stata uses <code>ivreg2</code> (Kleibergen-Paap) <code>widstat</code> F-statistic	
Error Attribution	Paper underspecifies code by not disclosing how party is coded in data		Paper underspecifies code by not disclosing how party is coded in data		Paper underspecifies F-statistic type	

Table 3. Examples of discrepancies of paper underspecification. *The examples illustrate how agent behavior can vary with an underspecified task.*

to a discrepancy or be silently resolved, depending on the agent. Consequently, higher-performing agents tend to exhibit fewer human-attributed discrepancies, not because the underlying issues are absent, but because they are more often resolved correctly during reimplementation.

5.4 Robustness to repeated runs and pre-training leakage

This section evaluates two dimensions of robustness in the results. First, since reasoning-based agentic systems are stochastic, we evaluate stability of performance across multiple runs. Second, some papers may have been included in pre-training, potentially leading to memorization, and so we compare performance for a selection of papers published before and after the model knowledge cutoffs.

Re-run stability. Because agentic systems are stochastic, we assess how sensitive reproduction outcomes are to repeated runs of the same task. We rerun the full pipeline twice for 20 randomly selected papers using Claude Code and Codex GPT-5.4, yielding three runs per paper (Figure 8). Overall, results are stable: more than 80% of tables exhibit a grade spread of at most one across runs. Exact agreement is common, and large deviations are rare, indicating that most reproduction outcomes are robust to randomness in the agent’s execution.

At the same time, stability at the table level masks meaningful variation at the coefficient level. Figure 8b shows the distribution of within-agent-cell across-run changes in coefficient estimates, in units of the original paper’s standard error. This chart suggests significant noise within agent systems, as about half of estimated coefficients are statistically different from themselves (that is, the coefficient from the same table cell) across run pairs. This suggests that while overall conclusions are stable, fine-grained comparisons between agents should be interpreted with caution, especially when differences are small.

Leakage analysis. To assess the potential role of pre-training leakage—where models may have been exposed to results during training—we compare performance on papers published before and after the model knowledge cutoff. Specifically, we reproduce five *EJ* papers published before and five published after the cutoff using Claude Code (Opus 4.6) and Codex (GPT-5.4), and test for differences in performance between the two groups.

Figure 9 reports the results for mean table grades with bootstrapped confidence intervals. For both Claude and GPT, we find no statistically significant difference in performance between pre- and post-cutoff papers. If anything, post-cutoff performance is slightly higher. The absence of a pre/post difference suggests that pre-training leakage is unlikely to be the primary driver of performance.

Mechanistically, such leakage would also be difficult to exploit: any memorized results would need to be translated into executable code, and likely across programming languages, since the agents generate Python while the original replication packages use other languages. Nonetheless, these results should be interpreted with caution, as the sample size is small, and

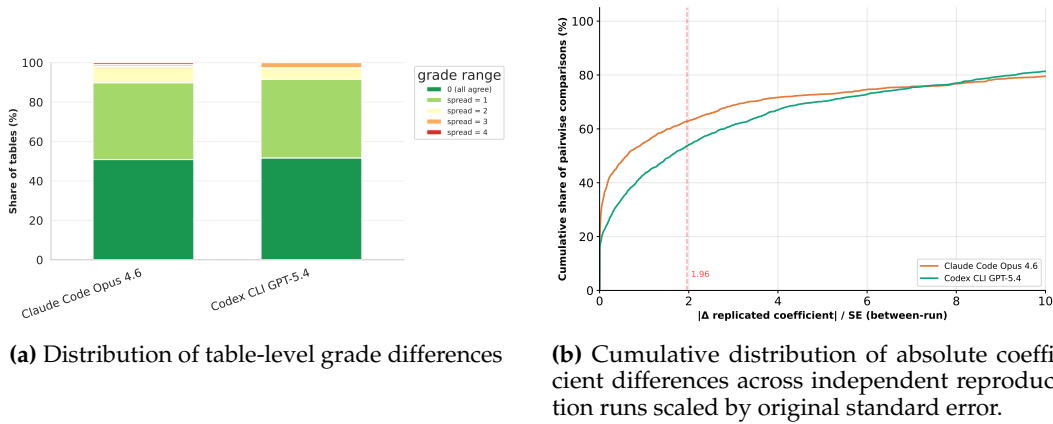


Figure 8. Stability of results between multiple reproduction runs. Both agents reproduce 20 random papers two additional times to a total of three runs. The grade range is computed as the difference between highest and lowest grade between runs excluding F grades. For the Panel B each coefficient cell that an agent replicated in two or more runs, the pairwise absolute difference between the replicated values, normalized by the original paper’s standard error is computed. The CDF reports the share of pairwise comparisons at or below each threshold. Additional figures can be found in Appendix A.5.

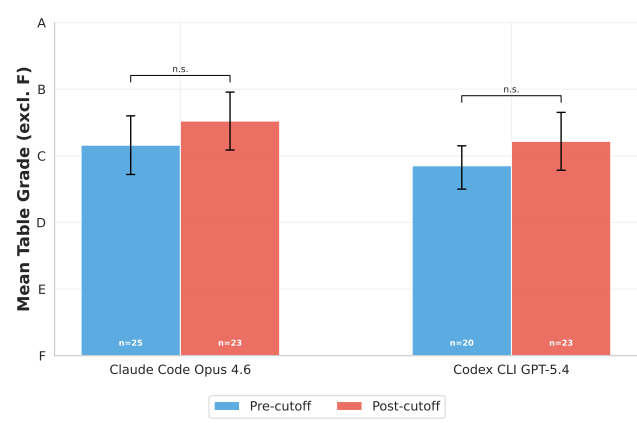


Figure 9. Pre-training leakage evaluation. Average table grades for sample of papers published before and after the model knowledge cutoff. No statistical difference suggests that performance of models in main analysis is not driven by pre-training leakage.

the selected papers have not undergone independent reproducibility verification. Moreover, this analysis does not rule out leakage arising during post-training (e.g., through fine-tuning or reinforcement learning on evaluation-relevant data). Still, overall, this analysis supports the interpretation that the observed performance reflects genuine reimplementations rather than retrieval of memorized outputs.

6 Conclusion

This paper establishes that AI agents can, in many cases, recover social science results from a paper’s text and data alone. Across a diverse set of empirical papers, agents are able to reconstruct a large share of published results without access to the original code, demonstrating substantial progress toward automated reproducibility. While this first generation of agents delivers imperfect results, the overall failure modes reveal a clear pattern. Going forward, the primary bottleneck may not be model capability, but the way social science methods are documented. In many cases, discrepancies arise because key implementation details are underspecified or omitted entirely, forcing agents to infer choices (just as a human researcher would under the same constraints). As a result, the ability to reproduce results reflects a combination of both the agent’s capabilities and the paper’s clarity and completeness.

We began this paper by articulating the common view that the scientific paper is the source of truth for scientific claims. However, many of the discrepancies we observe arise from missing or ambiguous specifications, such as variable coding choices, data filters, or estimation procedures. In this light, it is worth reconsidering what we mean by the “source of truth.” For the purpose of reproducing results, that role is already fulfilled by the code, which provides a complete, machine-readable specification of the analysis pipeline and resolves ambiguities that the paper often leaves implicit. This suggests a clearer division of roles: code serves as the authoritative representation of *what was done*, while the paper explains *why* those choices were made, defining variables, motivating identification strategies, and interpreting results in context.

This distinction has practical implications for how reproducibility is evaluated. Rather than expecting papers to fully encode implementation details, a more robust approach is to require explicit alignment between the narrative and the underlying code. Agentic systems provide a natural mechanism for this: they can translate papers into executable code, compare outputs, and surface inconsistencies between the two representations. Thus, automated reproduction can serve as both a tool for verification and a diagnostic, identifying where methods are underspecified, where implementations diverge, and where assumptions must be inferred. Aligning research practices with the requirements of both human and automated interpretation—through more explicit, structured, and complete method descriptions—may therefore play an important role in improving reproducibility at scale.

This perspective extends naturally to a broader spectrum of automated scientific tasks, each removing additional structure from the problem. Moving beyond reproduction with shared data raises a series of increasingly demanding questions. What if the data are unavailable, requiring agents to recover or reconstruct them? What if only the research question or hypothesis is given, and methods must be inferred rather than followed? What if agents are tasked not just with reproducing results, but with refining the analysis through specification checks, falsification of identification assumptions, or exploration of underlying mechanisms? And what if the goal shifts from reproduction to replication — asking questions or applying methods with new data?

While tackling these extensions could help address ongoing concerns about reproducibility in empirical research (e.g. Brodeur et al., 2024; 2026), such progress would also involve

pushing agentic systems further toward active participation in the research process. This shift requires us to define new criteria for scientific validity, including how to assess identification, robustness, and the reliability of conclusions generated without direct human oversight.

References

- Anthropic. Models overview. Anthropic API Documentation, 2026. URL <https://platform.claude.com/docs/en/about-claude/models/overview>. Accessed: 2026-04-01.
- Abel Brodeur and Bruno Barbarioli. The replication engine. Institute for Progress, 2025. URL <https://ifp.org/the-replication-engine/>.
- Abel Brodeur, Derek Mikola, Nikolai Cook, Thomas Brailey, Ryan Briggs, Alexandra de Gendre, Yannick Dupraz, Lenka Fiala, Jacopo Gabani, Romain Gauriot, Joanne Haddad, Ryan McWay, Joel Levin, Magnus Johannesson, Edward Miguel, Lennard Metson, Jonas Minet Kinge, Wenjie Tian, Timo Wochner, Sumit Mishra, Joseph Richardson, Giulian Etingin-Frati, Alexi Gugushvili, Jakub Procházka, Myra Mohnen, Jakob Möller, Rosalie Montambeault, Sébastien Montpetit, Jason Collins, Sigmond Ellingsrud, Alexander Kustov, Louis-Philippe Morin, Todd Morris, Erlend Fleisje, Elaheh Fatemi-Pour, Scott Moser, Matt Woerman, Tim Ölkers, Fabio Motoki, Anders Kjelsrud, Lucija Muehlenbachs, Andreea Musulan, Christian Czymara, Hooman Habibnia, Alexander Coppock, Idil Tanrisever, Marco Musumeci, Nicholas Rivers, and Rachel Joy Forshaw. Mass reproducibility and replicability: A new hope. Technical Report 107, Institute for Replication (I4R), apr 2024. URL <https://hdl.handle.net/10419/289437>.
- Abel Brodeur, Derek Mikola, Nikolai Cook, Lenka Fiala, Thomas Brailey, Ryan Briggs, Alexandra de Gendre, Yannick Dupraz, Jacopo Gabani, Romain Gauriot, Joanne Haddad, Goncalo Lima, Jörg Ankel-Peters, Anna Dreber, Douglas Campbell, Lamis Kattan, Diego Marino Fages, Fabian Mierisch, Pu Sun, Taylor Wright, Marie Connolly, Fernando Hoces de la Guardia, Magnus Johannesson, Edward Miguel, Lars Vilhuber, Alejandro Abarca, Mahesh Acharya, Sossou Simplicie Adjisse, Ahwaz Akhtar, Eduardo Alberto Ramirez Lizardi, Sabina Albrecht, Synøve Nygaard Andersen, Zubaria Andlib, Falak Arora, Thomas Ash, Etienne Bacher, Sebastian Bachler, Félix Bacon, Manuel Bagues, Timea Balogh, Alisher Batmanov, Mara Barschkett, Barış Kaan Basdil, Jaromír Baxa, Sascha O Becker, Monica Beeder, Louis-Philippe Beland, Abdel-Hamid Bello, Daniel Benenson Markovits, Grant Benjamin, Thomas Bergeron, Moussa P Blimpo, Marco Binetti, Carl Bonander, Joseph Bonneau, Endre Borbáth, Nicolai Borgen, Solveig Topstad Borgen, Jonathan Borowsky, Elisa Brini, Myriam Brown, Martin Brun, Stephan Bruns, Nino Buliskeria, Andrea Calef, Alistair Cameron, Pamela Campa, Santiago Campos-Rodríguez, Giulio Giacomo Cantone, Fenella Carpena, Perry Jess Carter, Paul Castañeda Dower, Ondrej Castek, Jill Caviglia-Harris, Gabriella Chauca Strand, Shi Chen, Sya In Chzhen, Jong Chung, Jason Collins, Alexander Coppock, Hugo Cordeau, Ben Couillard, Jonathan Crechet, Lorenzo Crippa, Jing Cui, Christian Czymara, Haley Daarstad, Danh Chi Dao, Daniel Dao, Marco David Schmandt, Astrid de Linde, Lucas De Melo, Lachlan Deer, Micole De Vera, Velichka Dimitrova, Jan Fabian Dollbaum, Jan Matti Dollbaum, Michael Donnelly, Luu Duc Toan Huynh, Tsvetomira Dumbalska, Jamie Duncan, Kiet Tuan Duong, Thibaut Duprey, Christoph Dworschak, Sigmund Ellingsrud, Ali Elminejad, Yasmine Eissa, Andrea Erhart, Giulian Etingin-Frati, Elaheh Fatemipour, Alexa Federice, Jan Feld, Guidon Fenig, Mojtaba Firouzjaeiangalougah, Erlend Fleisje, Alexandre FortiFriter-Chouinard, Julia Francesca Engel, Nadjim Fréchet, Reid Fortier, Tilman Fries, Michael James Frith, Thomas Galipeau, Sebastian Gallegos, Areez Gangji, Xiaoying Gao, Cloé Garnache, Attila Gáspár, Evelina Gavrilova, Arijit Ghosh, Garreth Gibney, Grant Gibson, Geir Godager, Leonard Goff, Da Gong, Javier González, Jeremy D Gretton, Cristina Griffa, Idaliya Grigoryeva, Maja Grøtting, Eric Guntermann, Jiaqi Guo, Alexi Gugushvili, Hooman Habibnia, Sonja Häffner, Jonathan D Hall, Olle Hammar, Amund Hanson Kordt, Barry Hashimoto, Jonathan S Hartley, Carina I Hausladen, Tomáš Havránek, Harry He, Matthew Hepplewhite, Mario Herrera-Rodriguez, Felix Heuer, Anthony Heyes, Anson T Y Ho, Jonathan Holmes, Armando Holzknecht, Yu-Hsiang Dexter Hsu, Shiang-Hung Hu, Yu-Shiuan Huang, Mathias Huebener, Christoph Huber, Kim P Huynh, Zuzana

Irsova, Ozan Isler, Niklas Jakobsson, Raphaël Jananji, Tharaka A Jayalath, Michael Jetter, Jenny John, Rachel Joy Forshaw, Felipe Juan, Valon Kadriu, Sunny Karim, Edmund Kelly, Duy Khanh Hoang Dang, Tazia Khushboo, Jin Kim, Gustav Kjellsson, Anders Kjelsrud, Andreas Kotsadam, Jori Korpershoek, Lewis Krashinsky, Suranjana Kundu, Alexander Kustov, Nurlan Lalayev, Audrée Langlois, Jill Laufer, Blake Lee-Whiting, Andreas Leibing, Gabriel Lenz, Joel Levin, Peng Li, Tongzhe Li, Yuchen Lin, Ariel Listo, Dan Liu, Xuewen Lu, Elvina Lukmanova, Alex Luscombe, Lester R Lusher, Ke Lyu, Hai Ma, Nicolas Mäder, Clifton Makate, Alice Malmberg, Adit Maitra, Marco Mandas, Jan Marcus, Shushanik Margaryan, Lili Márk, Andres Martignano, Abigail Marsh, Isabella Masetto, Anthony McCanny, Emma McManus, Ryan McWay, Lennard Metson, Jonas Minet Kinge, Sumit Mishra, Myra Mohnen, Jakob Moeller, Rosalie Montambeault, Sébastien Montpetit, Louis-Philippe Morin, Todd Morris, Scott Moser, Fabio Yoshio Suguri Motoki, Lucija Muehlenbachs, Andreea Musulan, Marco Musumeci, Munirul Nabin, Karim Nchare, Florian Neubauer, Quan M P Nguyen, Tuan Nguyen, Viet Nguyen-Tien, Ali Niazi, Giorgi Nikolaishvili, Ardyn Nordstrom, Patrick Nüss, Angela Odermatt, Matt Olson, Henning Øien, Tim Ölkens, Miquel Oliver I Vert, Emre Oral, Christian Oswald, Ali Ousman, Ömer Özak, Shubham Pandey, Alexandre Pavlov, Martino Pelli, Romeo Penheiro, Ryungyung Park, Eva Pérez Martel, Tereza Petrovičová, Linh Phan, Alexa Prettyman, Jakub Procházka, Aqila Putri, Julian Quandt, Kangyu Qiu, Loan Quynh Thi Nguyen, Andaleeb Rahman, Carson H Rea, Adam Reiremo, Laëtitia Renée, Joseph Richardson, Nicholas Rivers, Bruno Rodrigues, William Roelofs, Tobias Roemer, Ole Rogeberg, Julian Rose, Andrew Roskos-Ewoldsen, Paul Rosmer, Barbara Sabada, Soodeh Saberian, Nicolas Salamanca, Georg Sator, Daniel Scates, Elmar Schlüter, Cameron Sells, Sharmi Sen, Ritika Sethi, Anna Shcherbiak, Moyosore Sogaolu, Matt Soosalu, Erik Ø Sørensen, Manali Sovani, Noah Spencer, Stefan Staubli, Renske Stans, Anya Stewart, Felix Stips, Kieran Stockley, Stephenson Strobel, Ethan Struby, John P Tang, Idil Tanrisever, Thomas Tao Yang, Ipek Tastan, Dejan Tatić, Benjamin Tatlow, Féraud Tchuisseu Seuyong, Rémi Thériault, Vincent Thivierge, Wenjie Tian, Filip-Mihai Toma, Maddalena Totarelli, Van-Anh Tran, Hung Truong, Nikita Tsoy, Kerem Tuzcuoglu, Diego Ubfal, Laura Villalobos, Julian Walterskirchen, Joseph Tao-Yi Wang, Vasudha Wattal, Matthew D Webb, Bryan S Weber, Reinhard Weisser, Wei-Chien Weng, Christian Westheide, Kimberly White, Jacob Winter, Timo Wochner, Matt Woerman, Jared Wong, Ritchie Woodard, Marcin Wroński, Myra Yazbeck, Gustav Chung Yang, Luther Yap, Kareman Yassin, Hao Ye, Jin Young Yoon, Chris Yurris, Tahreen Zahra, Mirela Zaneva, Aline Zayat, Jonathan Zhang, Ziwei Zhao, and Yaolang Zhong. Reproducibility and robustness of economics and political science research. *Nature*, 652(8108):151–156, April 2026.

Clayton Cohn, Siyuan Guo, Surya Rayala, Hanchen David Wang, Naveeduddin Mohammed, Umesh Timalisina, Shruti Jain, Angela Eeds, Menton Deweese, Pamela J. Osborn Popp, Rebekah Stanton, Shakeera Walker, Meiyi Ma, and Gautam Biswas. Evidence-decision-feedback: Theory-driven adaptive scaffolding for llm agents, 2026. URL <https://arxiv.org/abs/2602.01415>.

Anna Dreber and Magnus Johannesson. A framework for evaluating reproducibility and replicability in economics. *Economic Inquiry*, 63(2):338–356, 2025. doi: <https://doi.org/10.1111/ecin.13244>. URL <https://onlinelibrary.wiley.com/doi/abs/10.1111/ecin.13244>.

Nic Fishman and Gabriel Sekeres. Editorial screening when science is cheap, 2025. URL https://nfw.fish/static/papers/agentlic_specification.pdf.

Shubham Gandhi, Dhruv Shah, Manasi Patwardhan, Lovekesh Vig, and Gautam Shroff. Researchcodeagent: An llm multi-agent system for automated codification of research

-
- methodologies. In Qingyun Wang, Wenpeng Yin, Abhishek Aich, Yumin Suh, and Kuan-Chuan Peng (eds.), *AI for Research and Scalable, Efficient Systems*, pp. 3–37, Singapore, 2025. Springer Nature Singapore. ISBN 978-981-96-8912-5.
- Yifan Hou, Buse Giledereli, Yilei Tu, and Mrinmaya Sachan. Do vision-language models really understand visual language? In *Proceedings of the 42nd International Conference on Machine Learning, ICML'25*. JMLR.org, 2026.
- Chuxuan Hu, Liyun Zhang, Yeji Lim, Aum Wadhvani, Austin Peters, and Daniel Kang. REPRO-bench: Can agentic AI systems assess the reproducibility of social science research? In Wanxiang Che, Joyce Nabende, Ekaterina Shutova, and Mohammad Taher Pilehvar (eds.), *Findings of the Association for Computational Linguistics: ACL 2025*, pp. 23616–23626, Vienna, Austria, July 2025. Association for Computational Linguistics. ISBN 979-8-89176-256-5. doi: 10.18653/v1/2025.findings-acl.1210. URL <https://aclanthology.org/2025.findings-acl.1210/>.
- Tianyu Hua, Harper Hua, Violet Xiang, Benjamin Klieger, Sang T. Truong, Weixin Liang, Fan-Yun Sun, and Nick Haber. ResearchCodeBench: Benchmarking LLMs on implementing novel machine learning research code, 2025. URL <https://arxiv.org/abs/2506.02314>.
- Qian Huang, Jian Vora, Percy Liang, and Jure Leskovec. MAgentBench: Evaluating language agents on machine learning experimentation, 2023. URL <https://arxiv.org/abs/2310.03302>.
- Institute for Replication. Meta database, version 1. <https://i4replication.org/reports/>, 2024. Accessed: 2026-03-31.
- John P. A. Ioannidis. Why most published research findings are false. *PLoS Medicine*, 2(8): e124, 2005. doi: 10.1371/journal.pmed.0020124.
- Haitao Li, Qian Dong, Junjie Chen, Huixue Su, Yujia Zhou, Qingyao Ai, Ziyi Ye, and Yiqun Liu. Llm-as-judges: a comprehensive survey on llm-based evaluation methods. *arXiv preprint arXiv:2412.05579*, 2024.
- Chris Lu, Cong Lu, Robert Tjarko Lange, Jakob Foerster, Jeff Clune, and David Ha. The AI scientist: Towards fully automated open-ended scientific discovery, 2024. URL <https://arxiv.org/abs/2408.06292>.
- Olivia Miske, Anna Lou Abatayo, Mason Daley, Mirka Dirzo, Nicholas Fox, Noah Haber, Krystal M Hahn, Melissa Kline Struhl, Brinna Mawhinney, Priya Silverstein, Theresa Stankov, Andrew H Tyner, Matúš Adamkovič, Shilaan Alzahawi, Saule Anafinova, Eli Awtrey, Erick Axxe, James Bailey, Bert N Bakker, Akshaya Balaji, Gabriel Banik, František Bartoš, Henk Berkman, Zachariah Berry, Felix S Bethke, Timothy F Brady, Nate Breznau, Sara Capitan, Tabaré Capitán, Laura Caquelin, Kent Jason Cheng, William J Chopik, Gwen-Jiro Clochard, Tom Coupé, Jamie Cummins, Elif Gizem Demirag Burak, Jianhua Duan, Kevin M Esterling, Thomas R Evans, Nathan Fiala, James Field, Victor Gay, Jing Geng, Johanna Gereke, Ilka Helene Gleibs, Amélie Gourdon-Kanhukamwe, Dmitry Grigoryev, Nicholas Gunby, Paul H P Hanel, Sanghyun Hong, Sean Dae Houlihan, Nick Huntington-Klein, Kamil Izydorczak, Kristin Jankowsky, Kai Jonas, Pavol Kačmár, Hansika Kapoor, Sebastian Karcher, Marta Kołczyńska, David Kretschmer, Ljiljana Lazarevic, Katelin E Leahy, Jessica C Lee, Christopher Limnios, An-Chiao Liu, John Wills Lloyd, Ruben Lopez-Nicolas, Nigel Mantou Lou, Richard E Lucas, Maximilian Maier, Daniel J Mallinson, Marcel Martončík, Michael C McCall, Nikita Mehta, Esteban Méndez, Johannes Michalak, Daniel C Molden, Faisal Mushtaq, Claudia Neuendorf, Austin Lee

-
- Nichols, Gustav Nilsonne, Ernest O’Boyle, Jeewon Oh, Thomas Ostermann, Abiola Oye-banjo, Radoslaw Panczak, Yuri G Pavlov, Zoran Pavlović, Noemi Peter, Kim Peters, Nathaniel D Porter, Mariah Purol, Arathy Puthillam, Marco Ramljak, Arran T Reader, W Robert Reed, Jan Philipp Röer, Ivan Ropovik, Alexander O Savi, Kathleen Schmidt, Landon Schnabel, Eric L Sevigny, Samuel Shaki, Shishir Shakya, Andrew Soh, Angela Somo, Fatih Sonmez, Eirik Strømmland, Jordan W Suchow, Anna Szabelska, Anirudh Tagat, Melba Verra Tutor, Karolina Urbanska, Pieter Van Dessel, Elisabeth Julie Vargo, Diem Thi Hong Vo, Victor Volkman, Ke Wang, Aaron L Wichman, Jamal R Williams, Fabian Winter, Ferdinand Wintermantel, Nan Zhang, Ignazio Ziano, Cristina Zogmaister, Zorana Zupan, Brian A Nosek, and Timothy M Errington. Investigating the reproducibility of the social and behavioural sciences. *Nature*, 652(8108):126–134, April 2026.
- Bang Nguyen, Dominik Soós, Qian Ma, Rochana R. Obadage, Zack Ranjan, Sai Koneru, Timothy M. Errington, Shakhlo Nematova, Sarah Rajtmajer, Jian Wu, and Meng Jiang. ReplicatorBench: Benchmarking LLM agents for replicability in social and behavioral sciences, 2026. URL <https://arxiv.org/abs/2602.11354>.
- Brian A. Nosek, Tom E. Hardwicke, Hannah Moshontz, Aurélien Allard, Katherine S. Corker, Anna Dreber, Fiona Fidler, Joseph Hilgard, Melissa Kline Struhl, Michèle B. Nuijten, Julia M. Rohrer, Felipe Romero, Anne M. Scheel, Laura D. Scherer, Felix D. Schönbrodt, and Simine Vazire. Replicability, robustness, and reproducibility in psychological science. *Annual Review of Psychology*, 73:719–748, 2022. doi: 10.1146/annurev-psych-020821-114157.
- Brian A. Nosek et al. An open, large-scale, collaborative effort to estimate the reproducibility of psychological science. *Perspectives on Psychological Science*, 7(6):657–660, 2012. doi: 10.1177/1745691612462588.
- OpenAI. GPT-5.4 pro model. OpenAI API Documentation, 2026. URL <https://developers.openai.com/api/docs/models/gpt-5.4-pro>. Accessed: 2026-04-01.
- Samuel Schmidgall, Yusheng Su, Ze Wang, Ximeng Sun, Jialian Wu, Xiaodong Yu, Jiang Liu, Michael Moor, Zicheng Liu, and Emad Barsoum. Agent laboratory: Using LLM agents as research assistants. In Christos Christodoulopoulos, Tanmoy Chakraborty, Carolyn Rose, and Violet Peng (eds.), *Findings of the Association for Computational Linguistics: EMNLP 2025*, pp. 5977–6043, Suzhou, China, November 2025. Association for Computational Linguistics. ISBN 979-8-89176-335-7. doi: 10.18653/v1/2025.findings-emnlp.320. URL <https://aclanthology.org/2025.findings-emnlp.320/>.
- Minju Seo, Jinheon Baek, Seongyun Lee, and Sung Ju Hwang. Paper2Code: Automating code generation from scientific papers in machine learning. In *International Conference on Learning Representations (ICLR)*, 2026. URL <https://arxiv.org/abs/2504.17192>.
- Syed Mehtab Hussain Shah, Frank Hopfgartner, and Arnim Bleier. Automating computational reproducibility in social science: Comparing prompt-based and agent-based approaches, 2026. URL <https://arxiv.org/abs/2602.08561>.
- Chenglei Si, Diyi Yang, and Tatsunori Hashimoto. Can LLMs generate novel research ideas? A large-scale human study with 100+ NLP researchers. In *International Conference on Learning Representations (ICLR)*, 2025. URL <https://arxiv.org/abs/2409.04109>.
- Zachary S. Siegel, Sayash Kapoor, Nitya Nadgir, Benedikt Stroebel, and Arvind Narayanan. Core-bench: Fostering the credibility of published research through a computational reproducibility agent benchmark. *Transactions on Machine Learning Research (TMLR)*, 2024. URL <https://arxiv.org/abs/2409.11363>. Accepted.

-
- Giulio Starace, Oliver Jaffe, Dane Sherburn, James Aung, Jun Shern Chan, Leon Maksin, Rachel Dias, Evan Mays, Benjamin Kinsella, Wyatt Thompson, Johannes Heidecke, Amelia Glaese, and Tejal Patwardhan. PaperBench: Evaluating AI’s ability to replicate AI research, 2025. URL <https://arxiv.org/abs/2504.01848>.
- Maojun Sun, Ruijian Han, Binyan Jiang, Houduo Qi, Defeng Sun, Yancheng Yuan, and Jian Huang. A survey on large language model-based agents for statistics and data science. *The American Statistician*, 0(0):1–14, 2025. doi: 10.1080/00031305.2025.2561140. URL <https://doi.org/10.1080/00031305.2025.2561140>.
- Liyan Tang, Grace Kim, Xinyu Zhao, Thom Lake, Wenxuan Ding, Fangcong Yin, Prasann Singhal, Manya Wadhwa, Zeyu Leo Liu, Zayne Rea Sprague, Ramya Namuduri, Bodun Hu, Juan Diego Rodriguez, Puyuan Peng, and Greg Durrett. Chartmuseum: Testing visual reasoning capabilities of large vision-language models. In *Advances in Neural Information Processing Systems (NeurIPS) Datasets and Benchmarks Track*, 2025. URL <https://openreview.net/forum?id=V8pG2uIeI4>.
- GLM-5 Team, Aohan Zeng, Xin Lv, Zhenyu Hou, Zhengxiao Du, Qinkai Zheng, Bin Chen, Da Yin, Chendi Ge, Chengxing Xie, et al. Glm-5: from vibe coding to agentic engineering. *arXiv preprint arXiv:2602.15763*, 2026. URL <https://huggingface.co/papers/2602.15763>.
- Minyang Tian, Luyu Gao, Shizhuo Dylan Zhang, Xinan Chen, Cunwei Fan, Xuefei Guo, Roland Haas, Pan Ji, Kittithat Krongchon, Yao Li, Shengyan Liu, Di Luo, Yutao Ma, Hao Tong, Kha Trinh, Chenyu Tian, Zihan Wang, Bohao Wu, Yanyu Xiong, Shengzhu Yin, Minhui Zhu, Kilian Lieret, Yanxin Lu, Genglin Liu, Yufeng Du, Tianhua Tao, Ofir Press, Jamie Callan, Eliu Huerta, and Hao Peng. SciCode: A research coding benchmark curated by scientists, 2024. URL <https://arxiv.org/abs/2407.13168>.
- Andrew H Tyner, Anna Lou Abatayo, Mason Daley, Samuel Field, Nicholas Fox, Noah A Haber, Krystal M Hahn, Melissa Kline Struhl, Brinna Mawhinney, Olivia Miske, Priya Silverstein, Courtney K Soderberg, Theresa Stankov, Ahmed Abbasi, Christopher L Aberson, Balazs Aczel, Matúš Adamkovič, Nihan Albayrak, Peter J Allen, Michael Andreychik, Eli Awtrey, Erick Axxe, Flavio Azevedo, Miles D Bader, Bence Bago, James Bailey, Marjan Bakker, Gabriel Banik, George C Banks, Ernest Baskin, Anatolia Batruch, Annika Beatteay, Sophie M Behr, Nicholas Berente, Zachariah Berry, Jędrzej Białkowski, Bojana Bodroža, Laura Boeschoten, Miklos Bognar, Christian Bokhove, Diane Bonfiglio, Robin Bouwman, Timothy F Brady, Scott R Braithwaite, Gabriel Briceño Jiménez, Cameron Brick, Traci Bricka, Roman Briker, Annette N Brown, Gordon D A Brown, Robbie C M van Aert, Kathryn Caldwell, Sara Capitan, Tabaré Capitán, Jesse Chandler, Tessa Charles, Christopher R Chartier, Rahul Chawdhary, Kent Jason Cheng, William J Chopik, Bruce Clark, Victoria E Colvin, C Cozette Comer, Giulio Costantini, Tom Coupé, Jamie Cummins, Aneta Czernatowicz-Kukuczka, Joshua de Leeuw, David Dobolyi, James N Druckman, Jianhua Duan, Marin Dujmović, Daniel J Dunleavy, Patrick K Durkee, Cécile Emery, Kevin M Esterling, Thomas R Evans, Anna Fedor, Belén Fernández-Castilla, Nathan Fiala, James G Field, Nathan Fong, Miguel A Fonseca, Alexandra L J Freeman, Jeremy Freese, Sandra J Geiger, Jing Geng, Laura M Getz, Linda Marjoleine Geven, Ilka Helene Gleibs, Donna Pamella Gonzales, Janaki Gooty, Amélie Gourdon-Kanhukamwe, Cristina Greculescu, Siobhán M Griffin, Lusine Grigoryan, Martina Grunow, Nicholas Gunby, Braeden Hall, Paul H P Hanel, Erin E Hannon, Sam Harper, Marco Jürgen Held, Louis Hickman, Nathan C Higgins, Svenja Hippel, Sven Hoeppepner, Sanghyun Hong, Thomas J Hostler, Michael Inzlicht, Kamil Izydorczak, Bastian Jaeger, Kristin Jankowsky, Johannes Jarke-Neuert, Matthew Jensen, Biljana Jokić, Daniel Jolles, Phillip Jolly, Angela M Jones, Marie Juanchich, Pavol Kačmár, Hansika Kapoor, Andjela Keljanovic, Samjhana Koirala, Marta Kołczyńska, Dimitra Kouroupaki, Ulrich Kühnen, Michelangelo

Landgrave, Michael J Larson, Lyonel Laulié, Alice C E Lawrence, Joel M Le Forestier, Katelin E Leahy, Sungmok Lee, Jared Leslie, Savannah C Lewis, Christopher Limnios, Hause Lin, An-Chiao Liu, John Wills Lloyd, Elliot A Ludvig, Dermot Lynott, Jordan MacDonald, Peter Mallik, Daniel J Mallinson, Daniele Marinazzo, Corinna S Martarelli, Joshua Maticotta, Andrew McBride, Cillian McHugh, Gail McMillan, Esteban Méndez, Mitchell Metzger, Michalis P Michaelides, Johannes Michalak, Leticia Micheli, Jeremy K Miller, Marina Milyavskaya, Daniel C Molden, Ambar G Monjaras, David Moreau, Audrey Morrow, Cristóbal Moya, Liad Mudrik, Laetitia B Mulder, Katie A Munt, Arijit Nandi, Kathryn Nason, Carolin Nast, Gideon Nave, Heinrich H Nax, Florian Neubauer, Phuong Linh L Nguyen, Austin Lee Nichols, Gustav Nilsson, Ernest O’Boyle, Jule Oettinghaus, Jeewon Oh, Adoril Oshana, Thomas Ostermann, Rachel P Ostrowski, Abiola Oyebanjo, Radoslaw Panczak, Jamie Patrianakos, Ignacio Pavez, Yuri G Pavlov, Sofia Persson, Marco Perugini, Kim Peters, Constant Pieters, Vladimir Ponizovskiy, Nathaniel D Porter, Jason M Prenoveau, Danka Purić, Mariah F Purol, Arathy Puthillam, Kimberly A Quinn, Marco Ramljak, W Robert Reed, Michaela Ritchie, Margaret Ritzau, Sean Patrick Roche, Romina Rodela, Jan Philipp Röer, Ivan Ropovik, Jacob Rothschild, Justine Saal, Hani Safadi, Jason Samaha, Mary Sanchez, Soorya Sankaran, David Santos, Amanda C Sargent, Marian Sauter, Kathleen Schmidt, Landon Schnabel, Amber N Schroeder, Sebastian W Schuetz, Brendan A Schuetze, Michael Schulte-Mecklenbeck, Astrid Schütz, Eric L Sevigny, Ellie Shackleton, Richard M Shafraneck, Samuel Shaki, Shishir Shakya, Miroslav Sirota, Matthew Ryan Sisco, Maksim M Sitnikov, L Robert Slevc, Laura Smalarz, Colin Tucker Smith, Joel S Snyder, Nicolas Sommet, Fatih Sonmez, Barbara A Spellman, Natalia Stanulewicz-Buckley, George Stock, Chris N H Street, Eirik Strömeland, Tina Sundelin, Moin Syed, Anna Szabelska, Barnabas Szaszi, Ewa Szumowska, Anirudh Tagat, Susanne Täuber, Louis Tay, Stuti Thapa, Jason Thatcher, Domna Tsaklakidou, Lars Tummers, Elise Turkovich, Melba Verra Tutor, Karolina Urbanska, Anna Elisabeth van ’t Veer, Marcel van Assen, Niels van de Ven, Ruben van den Goorbergh, Elisabeth Julie Vargo, Leigh Ann Vaughn, Simine Vazire, Jentien M Vermeulen, Diem Thi Hong Vo, Victor Volkman, Eric-Jan Wagenmakers, Deliah Wagner, Lukasz Walasek, Frank Walter, Lara Warmelink, Liuqing Wei, Marie Isabelle Weißflog, Nicholas Weller, Aaron L Wichman, Jonathan Wilbiks, Jamal R Williams, Kelly Wolfe, Finnian Wort, Ryan Wright, Jesper N Wulff, Xindong Xue, Veronica X Yan, Yuzhi Yang, Sangsuk Yoon, Iris Žeželj, Yinxian Zhang, Ignazio Ziano, Cristina Zogmaister, Zorana Zupan, Rolf A Zwaan, Brian A Nosek, and Timothy M Errington. Investigating the replicability of the social and behavioural sciences. *Nature*, 652(8108):143–150, April 2026.

Zhen Wang, Fan Bai, Zhongyan Luo, Jinyan Su, Kaiser Sun, Xinle Yu, Jieyuan Liu, Kun Zhou, Claire Cardie, Mark Dredze, Eric P. Xing, and Zhiting Hu. FIRE-Bench: Evaluating agents on the rediscovery of scientific insights, 2026. URL <https://firebench.github.io/>.

Zirui Wang, Mengzhou Xia, Luxi He, Howard Chen, Yitao Liu, Richard Zhu, Kaiqu Liang, Xindi Wu, Haotian Liu, Sadhika Malladi, Alexis Chevalier, Sanjeev Arora, and Danqi Chen. Charxiv: Charting gaps in realistic chart understanding in multimodal LLMs. In *Advances in Neural Information Processing Systems (NeurIPS) Datasets and Benchmarks Track*, 2024. URL <https://openreview.net/forum?id=cy8mq7QYae>.

Yiqing Xu and Leo Yang Yang. Scaling reproducibility: An AI-assisted workflow for large-scale replication and reanalysis, 2026. URL <https://arxiv.org/abs/2602.16733>.

John Yang, Carlos E Jimenez, Alexander Wettig, Kilian Lieret, Shunyu Yao, Karthik R Narasimhan, and Ofir Press. SWE-agent: Agent-computer interfaces enable automated software engineering. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024. URL <https://arxiv.org/abs/2405.15793>.

Linhao Zhang, Tong Xia, Jinghua Piao, Lizhen Cui, and Yong Li. PaperRepro: Automated computational reproducibility assessment for social science papers, 2026. URL <https://arxiv.org/abs/2603.00058>.

Xuanle Zhao, Zilin Sang, Yuxuan Li, Qi Shi, Weilun Zhao, Shuo Wang, Duzhen Zhang, Xu Han, Zhiyuan Liu, and Maosong Sun. AutoReproduce: Automatic AI experiment reproduction with paper lineage, 2025. URL <https://arxiv.org/abs/2505.20662>.

A Main Appendix

A.1 Data overview

Step	N
I4R Replicate universe (unique papers)	109
Perfect computational reproduction	59
Sufficient or partial data availability	54
Has extractable tables	48
Final sample	48

Table A1. Sample Construction. Decomposition of papers from I4R.

Language	N papers	% papers	Total LOC	Mean LOC
Stata	26	54.2%	203,953	7,844
R	13	27.1%	24,462	1,881
MATLAB	1	2.1%	5,520	5,520
mixed	7	14.6%	21,637	3,091
All	48	100.0%	255,572	5,324

Table A2. Primary language of replication packages in the final sample. LOC is a line-count across all code files.

Table A3. Papers in Final Sample

Journal	Title	Year	Language	LOC
APSR	Flight to Safety: COVID-Induced Changes in the Intensity of Status Quo Preference and Voting Behavior	2021	R	2,372
APSR	Campaign Contributions and Roll-Call Voting in the U.S. House of Representatives: The Case of the Sugar Industry	2022	Stata	93
APSR	Can't We All Just Get Along? How Women MPs Can Ameliorate Affective Polarization in Western Publics	2022	R	716
JOP	Not All Elections Are Created Equal: Election Quality and Civil Conflict	2021	Stata	148
JOP	Antinormative Messaging, Group Cues, and the Nuclear Ban Treaty	2021	R	1,544
JOP	Changing Tides: Public Attitudes on Climate Migration	2021	R	3,720
JOP	What Makes Anticorruption Punishment Popular? Individual-Level Evidence from China	2021	R	1,328
JOP	Black Workers in White Places: Daytime Racial Diversity and White Public Opinion	2021	R	1,297
EJ	Non-Linearities, State-Dependent Prices and the Transmission Mechanism of Monetary Policy	2021	MATLAB	5,520
EJ	Understanding Ethnolinguistic Differences: The Roles of Geography and Trade	2021	Stata	1,362
EJ	Gender Differences in Cooperative Environments? Evidence from the U.S. Congress	2021	Stata	16,514
EJ	Pre-Colonial Warfare and Long-Run Development in India	2021	Stata	4,758
EJ	Spillover Effects of Intellectual Property Protection in the Interwar Aircraft Industry	2021	Stata	2,595
EJ	Why Don't Firms Hire Young Workers During Recessions?	2021	Stata	1,127
EJ	The Wheels of Change: Technology Adoption, Millwrights and the Persistence in Britain's Industrialisation	2022	Stata	1,240

Journal	Title	Year	Language	LOC
EJ	Peer Effects in Academic Research: Senders and Receivers	2022	mixed	5,606
	The Power of Hydroelectric Dams: Historical Evidence from the United States over the Twentieth Century	2022	Stata	7,598
QJE	War, Socialism, and the Rise of Fascism: an Empirical Exploration	2022	Stata	2,415
REStud	Who Chooses Commitment? Evidence and Welfare Implications	2021	mixed	4,823
	Exposure and Preferences: Evidence from Indian Slums	2020	R	2,399
AJPS	Re-Assessing Elite-Public Gaps in Political Behavior	2020	R	1,731
AJPS	Public Infrastructure and Economic Development: Evidence from Postal Systems	2021	mixed	1,074
AJPS	Hate Crimes and Gender Imbalances: Fears over Mate Competition and Violence against Refugees	2021	R	3,068
AJPS	Ascriptive Characteristics and Perceptions of Impropriety in the Rule of Law: Race, Gender, and Public Assessments of Whether Judges Can Be Impartial	2021	R	2,118
AJPS	Decentralization Can Increase Cooperation Among Public Officials	2021	mixed	1,028
AJPS	The Geography of Repression and Opposition to Autocracy	2021	Stata	8,345
AJPS	Talking Shops: The Effects of Caucus Discussion on Policy Coalitions	2021	R	890
AJPS	Parties as Disciplinarians: Charisma and Commitment Problems in Programmatic Campaigning	2021	mixed	1,445
AJPS	Indecent Disclosures: Anticorruption Reforms and Political Selection	2021	R	1,075
REStud	Immigration and Redistribution	2023	mixed	5,287
AER	Interaction, Stereotypes, and Performance: Evidence from South Africa	2022	Stata	12,428
AER	Vulnerability and Clientelism	2022	Stata	5,412
AER	Enabling or Limiting Cognitive Flexibility? Evidence of Demand for Moral Commitment	2023	Stata	4,656
AER	Market Access and Quality Upgrading: Evidence from Four Field Experiments	2022	Stata	3,828
AER	Evaluating Deliberative Competence: A Simple Method with an Application to Financial Choice	2022	Stata	3,178
AER	When a Doctor Falls from the Sky: The Impact of Easing Doctor Supply Constraints on Mortality	2023	Stata	5,230
AER	Can Technology Solve the Principal-Agent Problem? Evidence from China's War on Air Pollution	2022	Stata	1,864
AER: In-sights	Wage Cyclicalities and Labor Market Sorting	2022	Stata	16,880
AEJ: Applied	Assortative Matching at the Top of the Distribution: Evidence from the World's Most Exclusive Marriage Market	2022	Stata	12,957
AEJ: Applied	Historical Lynchings and the Contemporary Voting Behavior of Blacks	2022	Stata	1,432
	How Effective Are Monetary Incentives to Vote? Evidence from a Nationwide Policy	2021	Stata	7,523
AEJ: Macro	Declining Worker Turnover: The Role of Short-Duration Employment Spells	2021	mixed	2,374
AEJ: Policy	Multinationals' Sales and Profit Shifting in Tax Havens	2022	Stata	7,081
AEJ: Policy	School Spending and Student Outcomes: Evidence from Revenue Limit Elections in Wisconsin	2022	Stata	10,813
AEJ: Policy	The Long-Run Effects of Sports Club Vouchers for Primary School Children	2022	Stata	29,492
AEJ: Policy	How Do Beliefs about the Gender Wage Gap Affect the Demand for Public Policy?	2022	Stata	34,984

Journal	Title	Year	Language	LOC
AJPS	Entertaining Beliefs in Economic Mobility	2022	R	2,204

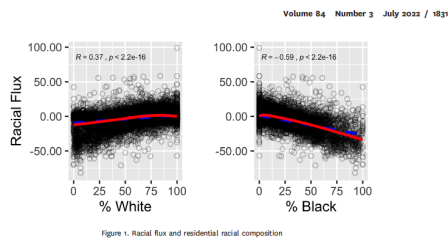
Step	Name	Inputs	Outputs
1	Extraction	Paper PDF Reproduction Package	Raw Data Extracted Methods Original Table Table Template
2	Reimplementation	Raw Data Extracted Methods Table Template	Script _i .py Table _i
3	Evaluation	Table _i Original Table	Comparison
4	Explanation	Paper PDF Reproduction Package Extracted Methods Script _i .py Comparison	Discrepancies & Error Sources

Table A4. Inputs and outputs for each step of the reproduction pipeline.

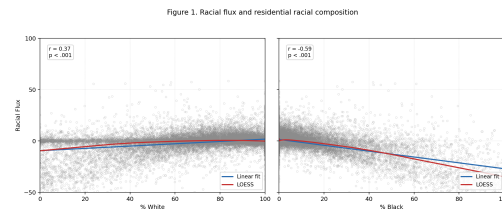
A.2 Reproduction of Figures

Our pipeline handles both tables and figures as reproduced outputs. For figures, the extractor generates a detailed description of the axes and plot structure without referencing any underlying data or results, and produces a simple Python template that recreates the plot and its styling from a data object (e.g., a dataframe). This design allows the reimplementer agent to focus on reconstructing the data while avoiding exposure to the original plot and minimizing the risk of result leakage, while maintaining stylistic consistency across runs.

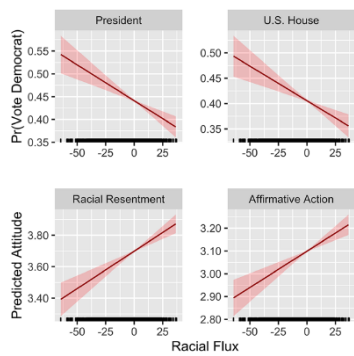
Evaluating reproduced figures is inherently more complex than comparing tables, as it requires assessing overall trends, individual data points, axes, and other visual features. Automating this process therefore relies on visual language models. However, whether such models can provide fair and consistent evaluations remains an open question. For this reason, we do not report quantitative results for figures. Figure A1 illustrates a representative example from GPT-5.4 in OpenCode, demonstrating the feasibility of the approach: the reproduced figures capture the main trends of the originals, though some visual discrepancies remain.



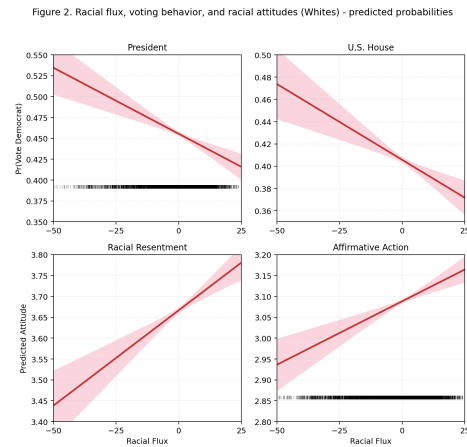
(a) Original Figure 1



(b) Reproduced Figure 1



(c) Original Figure 2



(d) Reproduced Figure 2

Figure A1. Agentic reproduction of figure results. Comparison between original and reproduced figures from the paper "Black Workers in White Places: Daytime Racial Diversity and White Public Opinion." The reproduction was conducted by OpenCode GPT-5.4.

A.3 Details on A-F Grading

Tables are evaluated through a comparison of aligned original and replicated outputs. First, cells are aligned across the two tables and replicated values may be rescaled when the comparison detects an approximate power-of-ten mismatch. Grades are assigned from the resulting cell-level differences. Because percentage differences are unstable when the original value is close to zero, the grading rule distinguishes between near-zero and non-near-zero cases. The grading rubric is outlined in Table A5

Grade	Near-zero original ($ x < 0.001$)	Otherwise
A	Absolute difference < 0.002 ; also assigned when both values are exactly zero	Percentage difference $< 2\%$
B	Absolute difference < 0.02	Percentage difference $< 20\%$
C	Absolute difference < 0.05	Percentage difference $< 40\%$
D	Absolute difference < 0.1	Percentage difference $< 60\%$
E	Absolute difference ≥ 0.1 , or signs differ	Percentage difference $\geq 60\%$, or signs differ
F	Assigned if either the original or replicated value is missing	

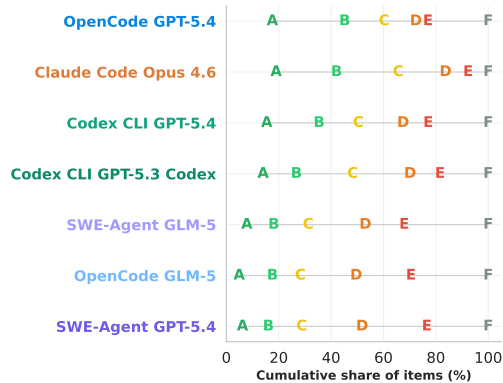
Table A5. Cell-level grading rules

Table-level grades are computed by converting cell grades to numeric values ($A = 5, B = 4, C = 3, D = 2, E = 1, F = 0$), averaging over non- F cells only, and then mapping the average back to a letter grade using fixed thresholds: $[4.5, 5] \rightarrow A, [3.5, 4.5) \rightarrow B, [2.5, 3.5) \rightarrow C, [1.5, 2.5) \rightarrow D$, and $[0.5, 1.5) \rightarrow E$. If all cells are graded F , the table grade is F . Paper-level grades are computed analogously, averaging across item grades using the same numeric mapping and the same thresholds, while excluding items marked unverifiable, items flagged as judge errors, and items graded F .

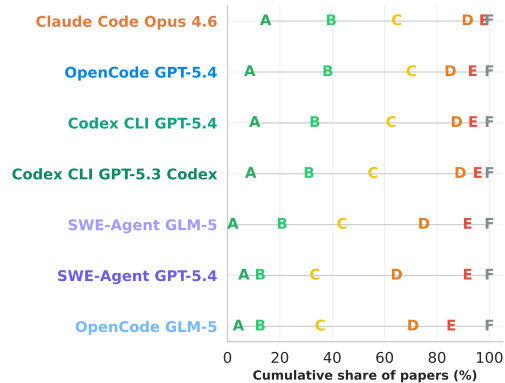
Figures are evaluated separately through a vision-capable language model and use the grading rubric presented in Figure A6.

Grade	Criterion
A	Visually indistinguishable
B	Patterns match, visible differences
C	Recognizable, noticeable gaps
D	Substantially different
E	Fundamentally different
F	Missing, blank, or not verifiable

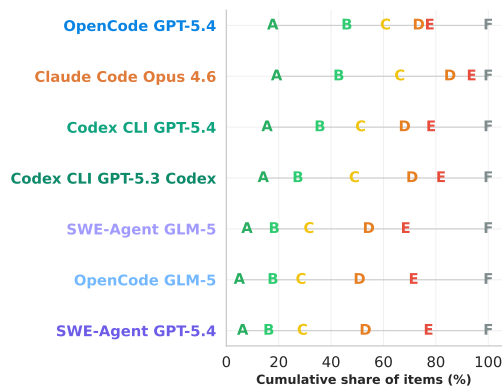
Table A6. Figure grading rubric



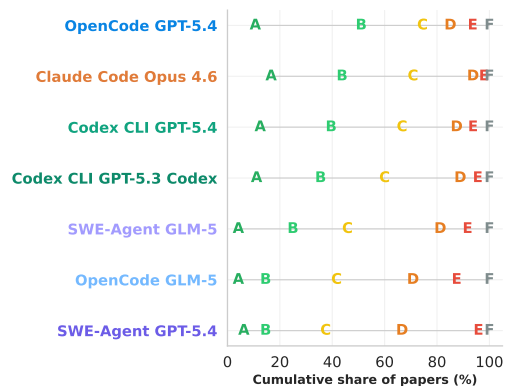
(a) Table-level grades including all F-graded cells.



(b) Paper-level grades including all F-graded tables.



(c) Table-level grades including F-grades for cells with at least one agent with a non-F result.



(d) Paper-level grades including F-grades for tables with at least one agent with a non-F result.

Figure A2. Table-level grades including F-grades. Aggregated performance naturally shrinks when including empty/not-reproduced tables/items. Notably the best-performing model/scaffold also switches on the paper level.

A.4 Additional Figures

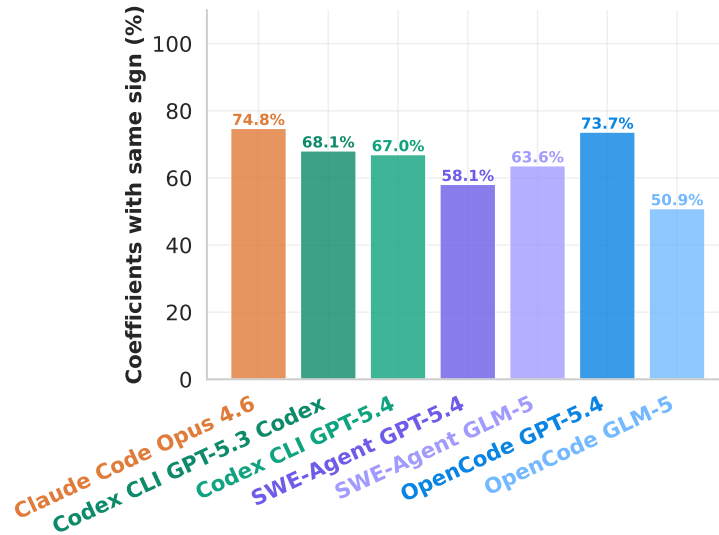


Figure A3. Share with correct sign, including missings. Share of reproduced coefficients with the same sign as the original paper including coefficients that are not being replicated in the denominator.

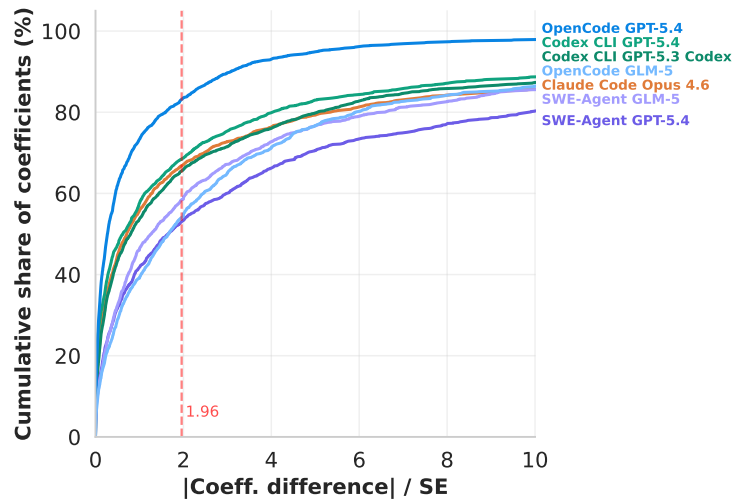


Figure A4. CDF of SE-adjusted coeff difference, excluding descriptive stats tables. Cumulative distribution of $|\text{coefficient difference}| / \text{SE}$ across approaches when removing descriptive tables. The dashed red line marks 1.96 (95% CI threshold. Neither the magnitude nor the order of results varies relevantly).

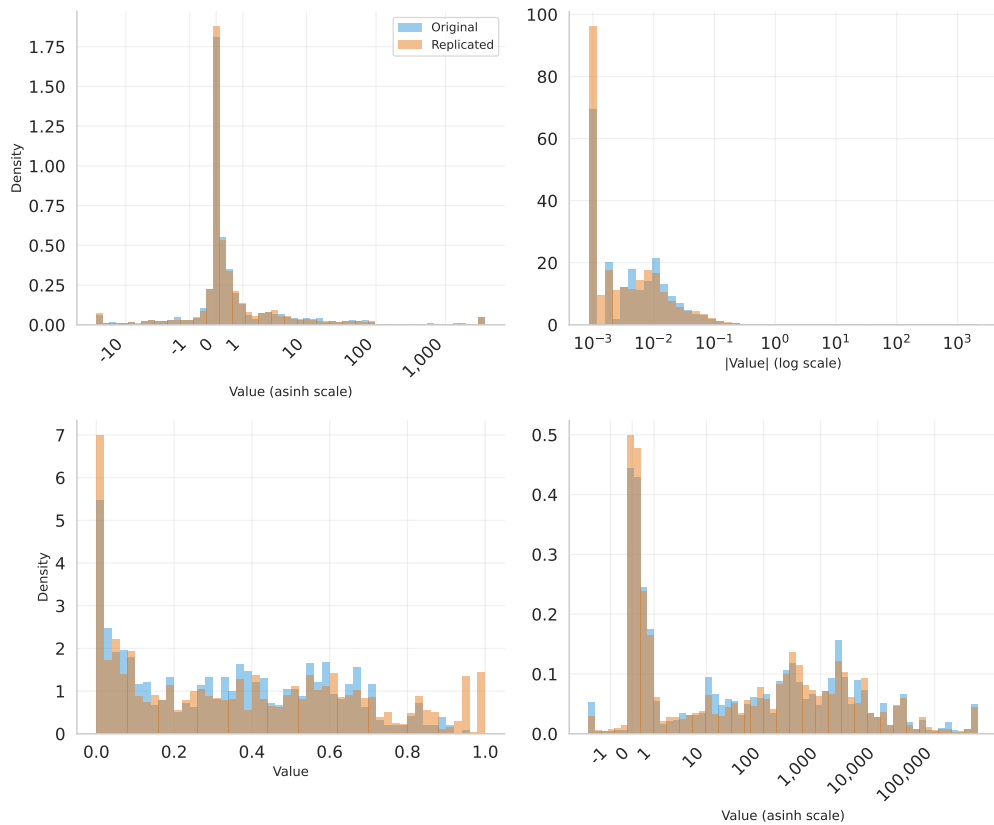


Figure A5. Histograms of original and reproduced statistics. Across various transformations of the statistics, the distributions of the original and reproduced statistics are quite similar.

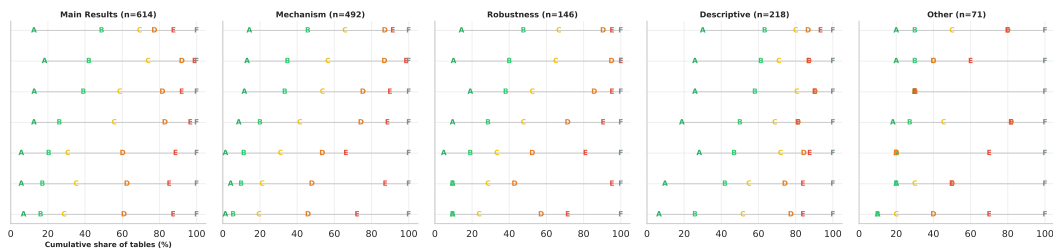


Figure A6. Grade distributions by table type. Including all F-grades.

A.5 Multi-run Stability

Based on 20 random papers re-run two times (three runs in total).

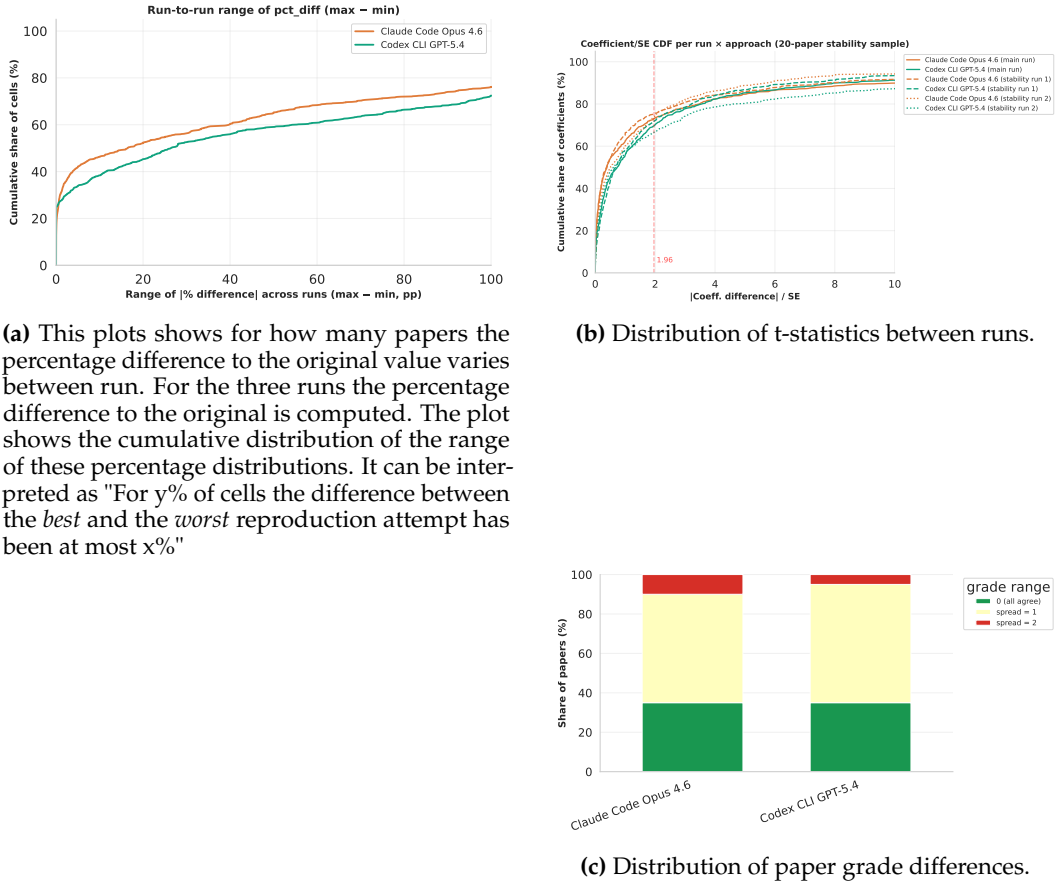


Figure A7. Additional statistics on stability between runs

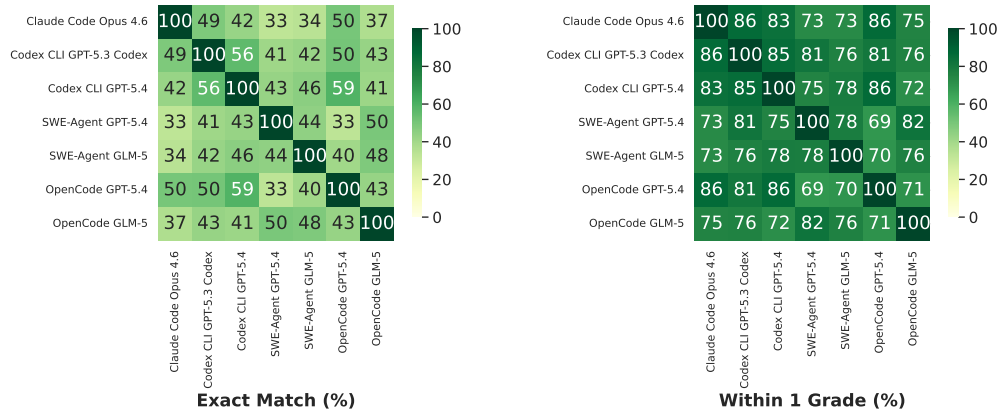


Figure A8. Inter-agent agreement at the paper level. Left: share of papers for which two agents assign the same grade. Right: share of papers for which grades are within one step.

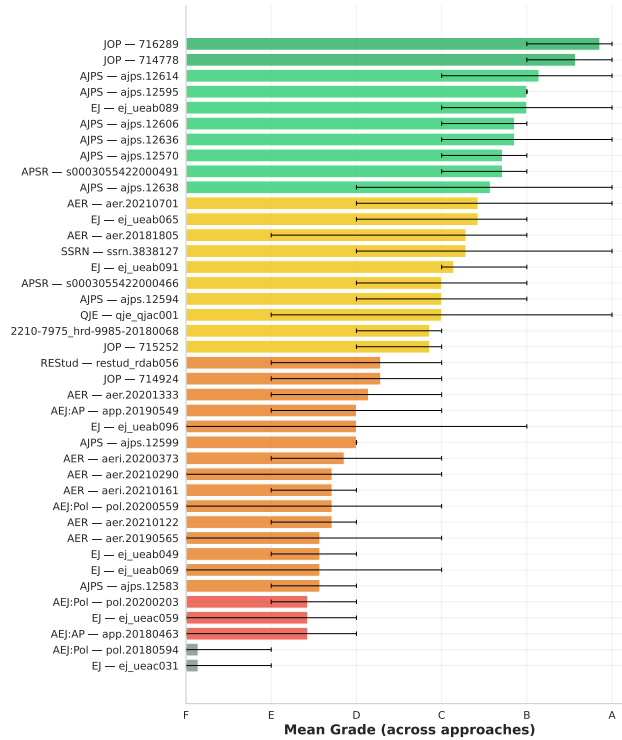


Figure A9. Within-paper performance range across all agent–model combinations, sorted by mean grade. Each bar spans the worst to best grade assigned by any agent.

A.6 Correlates of reproduction performance

We analyze potential sources of performance variation across papers and agents. Figure A8 shows how frequently two reproducing agents assign the same grade (left panel) or grades within one step of each other (right panel) for a given paper. This indicates that agents—particularly those developed by the same organization—tend to agree more than would be expected by chance, though between-agent performance nonetheless varies substantially.

Figure A9 offers a complementary perspective by displaying the within-paper performance range, which can also be interpreted as a measure of reproduction difficulty. The range

varies across papers but remains within two grade steps for the majority, highlighting that some papers are inherently harder to reproduce than others, while differences in agentic performance remain meaningful throughout.

Performance also varies across disciplines (Figure A10a): economics papers are generally reproduced less accurately than political science papers. However, since our sample is not random—having been selected by I4Replication based on study reproducibility—this pattern should not be interpreted as a disciplinary characteristic. A similar selection caveat applies to the comparison between performance and the original authors’ code language (Figure A10b): the finding that papers originally written in R tend to be reproduced more accurately than those in Stata or MATLAB should be interpreted with corresponding caution.

Finally, we examine whether overall dataset size or the volume of generated code are associated with reproduction performance. No systematic relationship is apparent between dataset size and grade, while a modest positive association is observed between the length of reproduced code and grade.

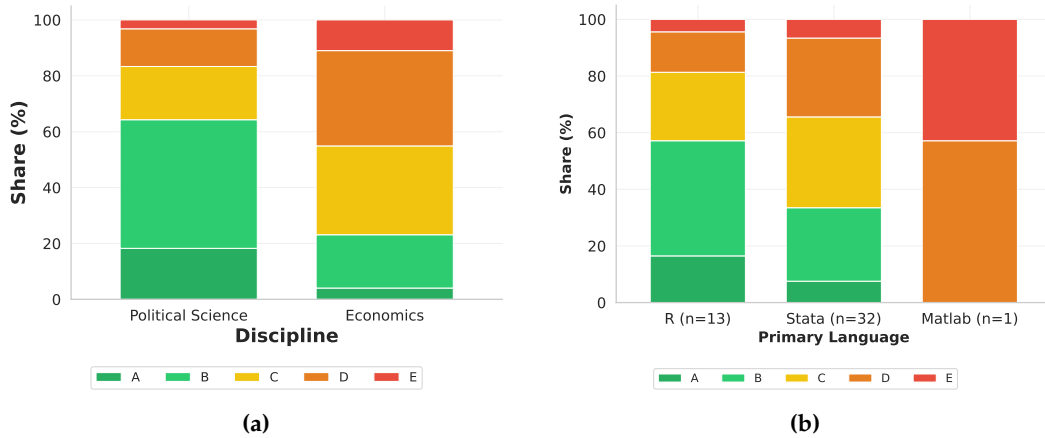


Figure A10. Paper-level grades by social science discipline and programming language used in the original reproduction package. Only papers with a single programming language considered.

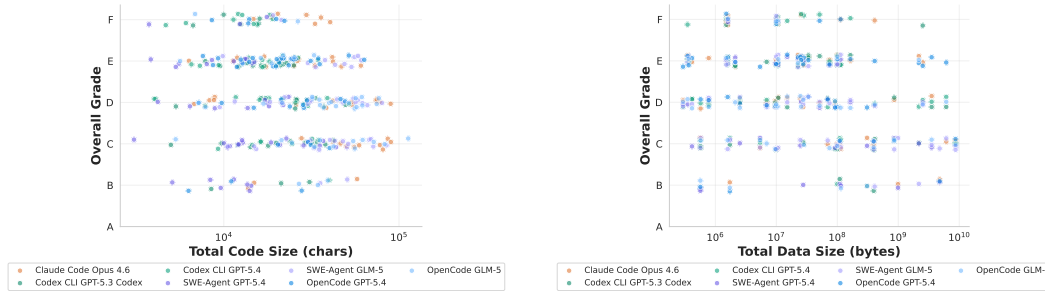


Figure A11. Paper-level grade by total code characters produced by the reproducing agent and the size of all datasets required for reproduction.

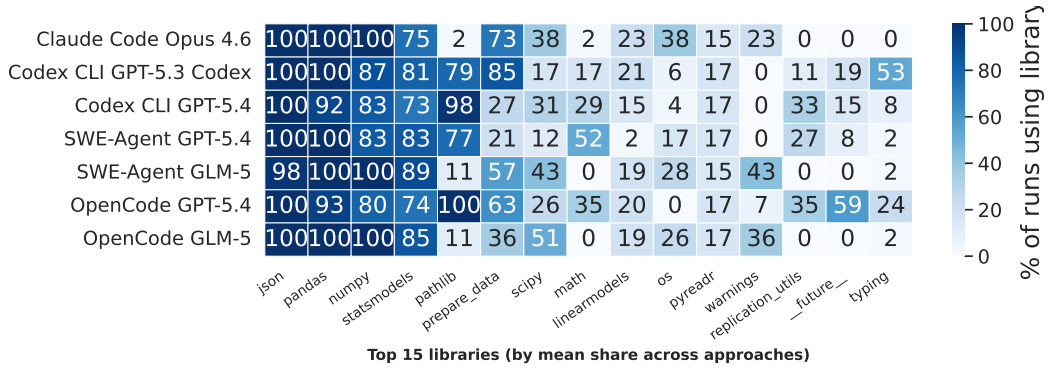


Figure A12. Incidence of Python imports across runs by agent. Note that *prepare_data* and *replication_utils* are python scripts that the agent created

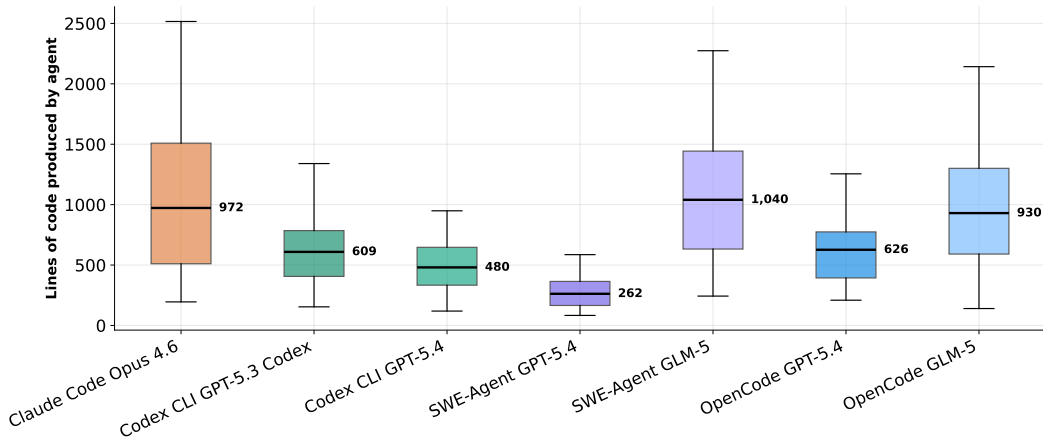


Figure A13. Distribution of lines of codes produced by agents. Box-and-whisker plots summarizing the number of lines of code produced by each agent system across runs.

A.7 Error source analysis

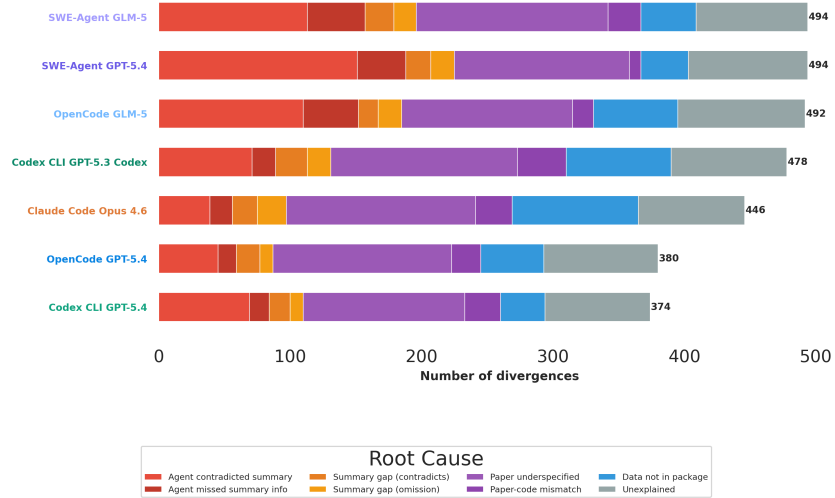


Figure A14. Comparison of error source attribution pipeline results between agents

Table A7. Example for output of error source pipeline identifying divergences, assigning severity, a stage of failure and a type of error for paper 10.1086_715252

Output	Sev.	Fail stage	Fail type	Agent	Description
Table 1	critical	Sum. \neq Agent	✗	Claude	Agent uses simple AMCE-difference approach instead of the parallel design causal mediation estimator (med.para) with covariates, and reports two-sided instead of one-sided p-values.
Table 1	critical	Code \neq Paper	○	Codex 5.3	The Python replicator substitutes the original parallel-design mediation/randomization-inference procedure with a pooled round-level product-of-coefficients OLS routine, so the indirect-effect estimates and their reported p-values are not computed the way the R code does.
Table 1	critical	Sum. \neq Agent	✗	Codex 5.4	The agent maps the mediator-specific natural and controlled columns to the wrong experiment indicators, so the natural columns use 'exp_id*_con' and the controlled columns use 'exp_id*_nat'.
Table 1	critical	Code \neq Paper	○	OC 5.4	The Python replicator replaces the original covariate-adjusted 'med.para.ri()' mediation analysis on 'exp_id_* == 1' versus '== 2' experiment arms with round-based treatment-control mean differences and separate within-round permutation tests.

continued on next page

(continued)

Output	Sev.	Fail stage	Fail type	Agent	Description
Table 1	critical	Code \neq Paper	○	SWE 5.4	The Python replicator does not implement the original mediation routine for Table 1 and instead fills every indirect-effect, parenthesized uncertainty, and total-effect cell from a simple punishment mean-difference plus bootstrap computed on hand-picked round subsets.
Table 1	medium	Code \neq Paper	○	OC GLM	The Python code reports HC1 standard errors in the inference rows, while the original R code reports randomization-inference tail probabilities ('indirect.greater' and 'total.greater').
Table 1	critical	Code \neq Paper	○	SWE GLM	The Python replicator does not implement the original 'med.para.ri' parallel-design mediation procedure for Table 1; instead it builds columns from ad hoc round filters and estimates indirect effects with two OLS regressions on 'rate_moral'/'rate_competence', omitting the original experiment-ID split, covariate set, and permutation-based inference.
Table 1	medium	Code \neq Paper	○	Codex 5.3	For total effects, the Python code estimates a simple treatment coefficient on pooled round subsets, while the R code estimates the total effect only on the treatment-only experiment arm selected by each 'exp_id_*' variable and includes the specified conjoint covariates.
Table 1	critical	Code \neq Paper	○	Codex 5.4	The agent does not implement the paper's mediation estimator with covariate adjustment and randomization-inference p-values; it instead fits a pooled 'vote ~ punishment * group' regression and reports cluster-robust normal p-values.
Table 1	medium	Code \neq Paper	○	OC GLM	The Python replicator replaces the mediation estimator with a simple difference between treatment coefficients from experiment arms, omitting the covariate-adjusted total-effect regression and the mediator-level weighted direct-effect calculation used in the R code.

Notes: Fail type: ✗ contradicts (direct conflict); ○ omission (upstream specifies, downstream silent); ? unclear. Agents: Claude = Claude Code Opus 4.6; Codex 5.3/5.4 = Codex CLI; OC = OpenCode; SWE = SWE-Agent; GLM = GLM-5.

A.8 Code snippets from reproduction agents

```
def main():
    df = main_sample(load_data())
    table = load_template(2)
    formulas = [
        "OUTCOME = tbul3_3rd + tcb + treat",
        "OUTCOME = treat + C(year_3rd) + C(bul3_3rd)",
        "OUTCOME = treat + C(year_3rd) + C(bul3_3rd) + C(cityno)",
    ]
    for col_idx, formula in enumerate(formulas):
        for outcome, row_coef, row_se in OUTCOMES:
            used = df[[outcome, "tbul3_3rd", "tcb", "treat", "year_3rd", "bul3_3rd", "cityno"]].dropna().copy()
            result = run_ols(formula, replace="OUTCOME", outcome, used)
            fill_coef_set(
                table,
                row_coef,
                row_se,
                col_idx,
                result.params["treat"],
                result.bs["treat"],
                result.pvalues["treat"],
            )
        set_cell(table, 28, col_idx, format_observations(len(used)), len(used))
    write_table(table, PATH("table_2.json"))
```

(a) Codex 5.4 — *table_2.py*

```
def run_ols(formula, data, cluster_col="cityno", weights=None):
    data = data.dropna().copy()
    if weights is None:
        res = smf.ols(formula, data=data).fit(
            cov_type="cluster",
            cov_kwds={"groups": data[cluster_col]},
        )
    else:
        res = smf.wls(formula, data=data, weights=data[weights]).fit(
            cov_type="cluster",
            cov_kwds={"groups": data[cluster_col]},
        )
    return res
```

(b) Codex 5.4 — *replication_utils.py*

```
def add_entropy_weights(aggates: pd.DataFrame, features: list[str]) -> pd.DataFrame:
    pre_period =
    aggates.loc[aggates["year_3rd"].isin([2006, 2007])]
    .groupby("municipality_id") [features]
    .mean()
    municipality_state = aggates[["municipality_id", "state_3rd"]].drop_duplicates().set_index("municipality_id")
    treated_ids = municipality_state.index[municipality_state["state_3rd"].eq(13)]
    control_ids = municipality_state.index[~municipality_state["state_3rd"].eq(13)]
    target_means = pre_period.loc[treated_ids, features].mean(axis=0)
    control_means = entropy_balance(pre_period.loc[control_ids, features], target_means)
    weights = pd.Series(1.0, index=pre_period.index, name="base_weight")
    weights.loc[control_ids] = control_weights
    return aggates.merge(weights, left_on="municipality_id", right_index=True, how="left")

def main() -> None:
    df = load_data()
    aggates = table7.aggates(df)
    table = load_json(TEMPLATE_PATH)
    for col_index, column_label in enumerate(table["column_labels"]):
        weighted = add_entropy_weights(aggates, VERSION_FEATURES[column_label])
        for outcome, coef_row, se_row in OUTCOME_ROWS:
            result = fit_ols(
                weighted,
                f"{outcome} = voucher + C(year_3rd) + C(municipality_id)",
                weight_col="base_weight",
            )
            bundle = coef_bundle(result, "voucher", add_stars=False)
            set_cell(
                table,
                coef_row,
                col_index,
                raw_text_bundle("coef_text"),
                numeric_value_bundle("coef"),
            )
        set_cell(
            table,
            se_row,
            col_index,
            raw_text_bundle("se_text"),
            numeric_value_bundle("se"),
        )
```

(c) OpenCode GPT 5.4 — *table_7.py*

```
def entropy_balance(
    control_features: pd.DataFrame,
    treated_target: pd.Series,
    *,
    maxiter: int = 10000,
    tol: float = 1e-10,
) -> np.ndarray:
    x = np.asarray(control_features, dtype=float)
    target = np.asarray(treated_target, dtype=float)

    def moment_gap(lmbda: np.ndarray) -> np.ndarray:
        scores = x @ lmbda
        scores -= scores.max()
        weights = np.exp(scores)
        weights /= weights.sum()
        return weights @ x - target

    sol = optimize.root(moment_gap, np.zeros(x.shape[1]), method="hybr", tol=tol)
    if not sol.success:
        raise RuntimeError(f"Entropy balancing did not converge: {sol.message}")
    scores = x @ sol.x
    scores -= scores.max()
    weights = np.exp(scores)
    weights /= weights.sum()
    weights *= len(weights)
    return weights
```

(d) OpenCode GPT 5.4 — *replication_utils.py*

Figure A15. Example screenshots from agents reproducing standard econometric methods. All screenshots are from reproduction attempts of the paper "The Long-Run Effects of Sports Club Vouchers for Primary School Children" (Marcus, Siedler & Ziebarth - AEJ Economic Policy (2022)). The screenshots show segments of the reproduced code for this paper. Panel (a) shows the main regression using a TWFE model for DiD analysis. Functions like *run_ols*, *format_observation*, *write_table* are created as helper functions in a separate file. Panel (b) shows such a function to be used in multiple table creation files. Panel (c) shows the method of an alternative empirical approach via an entropy-balanced synthetic control, applying a shared function shown in Panel (d).

B Trace analysis and guardrails

B.1 Setup

Each agent operates in a workspace containing: `TASK.md` (replication instructions extracted from the paper), `methodology_summary.json` (structured methodology context), `table_templates/*.json` (expected output structure), a symlink to the original dataset, and agent-specific instruction files. These files contain paper-derived content by design. The original replication package and paper PDF reside on the same machine but outside the workspace boundary. Web access was restricted using the agent scaffold’s configuration to the extent possible.

B.2 Trace analysis

Runs are analyzed by parsing execution logs and extracting tool calls. The agents differ in how they expose file operations: Claude Code and OpenCode provide structured Read/Write/Edit tools alongside bash, Codex exposes bash together with a programmatic file editor, and SWE-Agent operates primarily through bash.

To normalize across these interfaces, each bash command is split into sub-commands at common boundaries (e.g., `&&` chaining), while preserving heredoc bodies as atomic units. Each sub-command is then classified by its first token into one of:

- `exec`: executing code (e.g., `python`, `Rscript`, `make`)
- `read`: reading files (e.g., `cat`, `head`, `sed`)
- `navigation`: filesystem navigation (e.g., `ls`, `cd`, `pwd`)
- `search`: searching (e.g., `grep`, `find`, `rg`)
- `write`: writing files (e.g., `echo`, `cp`, `tee`, `cat` with `>`)
- `other`: everything else

For each run we measure two quantities: the *action count*, i.e. the number of tool calls (or sub-commands in the case of bash), and the *tool-call character count*, i.e. the number of characters the LLM generated per tool call as illustrated in Figure 6 in the main body.

B.3 Guardrail audit

For the replication runs we extract the event trace (all tool calls with inputs and truncated outputs) and the workspace file tree and contents. Guardrail adherence is assessed via two independent signals:

Regex scan. We extract all absolute file paths and URLs from the raw event trace using regular expressions. Each path is classified as: *allowed_workspace* (inside the workspace), *allowed_data* (the mounted dataset symlink), *forbidden_replication_package* (inside `papers/<id>/replication_package/`), *forbidden_paper_pdf* (the paper file), or *forbidden_external* (any other path outside the workspace). URLs and web-access keywords (`curl`, `wget`, `requests`) are flagged separately. This provides a deterministic baseline independent of LLM judgment.

LLM review. In an effort to flag additional runs for potential violations we prompt GPT-5.4-mini with the normalized event trace and workspace artifacts to identify access breaches. Each breach is classified as:

- *external_result_lookup*: fetching results from the web, APIs, or forbidden files outside the workspace;
- *forbidden_paper_access*: reading the published paper or prior replication results;
- *forbidden_code_access*: reading original replication scripts or code from outside the workspace.

Each breach includes the artifact and line range, severity (low/medium/high), confidence (low/medium/high), and an evidence chain. The overall run assessment ranges from *clean* to *severe_violation*,

with *insufficient_evidence* for runs where logs are missing or truncated. The model is explicitly instructed that all workspace files are allowed inputs and that incorrect replications are not breaches. In practice we find the LLM based review to be over-inclusive; severe suspected cases are manually inspected and rerun if breaches were found.

B.4 Hardcoding audit

For each run, we provide GPT-5.4-mini with all generated Python scripts from the workspace alongside the methodology summary and table templates for context. The model identifies numeric literals that appear in the script’s output as statistical results but have no computation path from the dataset. Each run is classified from *clean* to *severe_hardcoding*, with *insufficient_evidence* when Python files are missing or incomplete. In Figure A16 we classify the runs per agent scaffolding and model into their respective hardcoding categories. To assess whether hardcoding inflates replication grades, we compare runs classified as clean ($n = 299$) against all non-clean runs ($n = 37$). Clean runs achieve a mean grade of 3.02 ($\approx C$) compared to 2.50 ($\approx D$) for non-clean runs. Runs flagged for hardcoding perform no better—and if anything worse—than clean runs, providing no evidence that hardcoding inflates replication outcomes, though we acknowledge that hardcoding remains an undesirable behavior.

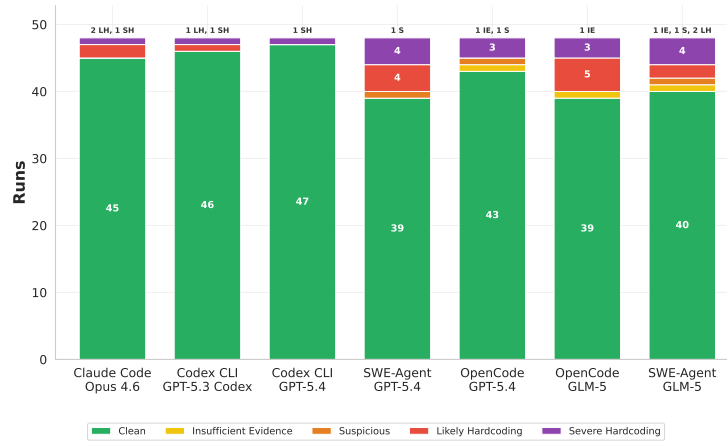


Figure A16. Hardcoding assessment outcomes by model and agent scaffold. Each stacked bar shows the number of runs classified from *clean* to *severe hardcoding*.

C Prompts

C1. Extraction Prompt

You are a methodological extraction specialist for social science research papers.

Your task: extract the methodology, structure, and specifications from academic papers so that a replicator who cannot see the paper can reproduce every table and figure. Extract NO results.

1. What to extract

- Research questions
- Data description, source, sample size, time period
- Data processing steps: BOTH general steps (applied to all analyses) AND per-table/per-figure steps (sample restrictions, variable construction, or filtering specific to one table or figure)
- Per-table: exact column headers, exact row labels, panel structure, caption, notes, regression specifications, data source (if different tables use different datasets)
- Per-figure: full caption, plot type, axis labels, legend/series entries, subplot structure, visual details (approximate axis ranges, reference lines, line styles, color conventions)
- Per-regression-spec: model type, dependent variable, independent variables, controls, fixed effects, clustering, sample restrictions, equation, variable definitions

Extract ALL tables and ALL data-based figures. Skip only purely conceptual visualizations (flow diagrams, frameworks, screenshots, photos, maps). Do not skip tables because they seem simple - summary statistics, balance tables, and cross-tabulations must all be extracted.

2. What NOT to extract (results)

Never include: regression coefficients, standard errors, t-statistics, p-values, significance stars, point estimates, confidence intervals, effect sizes, or descriptions of direction/magnitude.

DO include: table/figure structure (headers, labels, panels), design parameters (sample sizes, thresholds, variable names), and descriptive labels ("N", "R-squared", "Controls", "Yes"/"No").

3. Data processing and cleaning steps

This is critical for replication. Extract every step needed to go from the raw data files to each analysis-ready sample. Separate these into:

****General steps**** (applied before any specific analysis):

- File loading, merging, and reshaping (which files, which keys)
- Variable construction (new columns derived from raw data)
- Data cleaning: missing value treatment, outlier removal, winsorization
- Sample filtering: time period restrictions, geographic restrictions, demographic restrictions
- Transformations: log transforms, standardization, encoding

****Per-table / per-figure steps**** (use ``sample_restrictions`` on each table/figure spec):

- Additional filtering specific to one table (e.g., "Table 3 restricts to males aged 25-64")
- Subsample definitions for different columns (e.g., "Columns 1-3 use the full sample, columns 4-6 restrict to the treatment group")
- Variable construction specific to one analysis (e.g., "For Figure 2, compute rolling 30-day averages of daily returns")

Use the actual variable names from the dataset where possible.

4. Regression specifications

Attach each regression spec to the table or figure that displays its results. Include one spec per distinct model variant (e.g., different controls or subsamples across columns).

For each spec:

- Use actual variable names from the paper, not generic descriptions.
Write ``dependent_var``: "log(wage)", not ``dependent_var``: "outcome variable".
- ``controls``: enumerate every control variable individually. Never write aggregate descriptions like "socio-demographic controls"--- list each variable: ["age", "female", "education_level", ...].
- ``equation_latex``: copy the exact equation from the paper in LaTeX notation, preserving all subscripts, Greek letters, and notation. If no explicit equation is given, write one from the prose description. Example:
$$\hat{Y}_i = \beta_0 + \beta_1 X_i + \gamma Z_i + \delta_j + \epsilon_i$$
- ``variable_definitions``: define every symbol verbally, separated by semicolons. Include all information the paper provides: definition, construction, coding, unit, sign convention, omitted categories, data source. Only include what is explicitly stated--- do not infer.

```

Example: `A_i^T`: acceptance of targeted tax (1 if not "No", as defined in Section 2.3);
c_i: income-threshold fixed effects (4 levels: bottom 20/30/40/50 percentile)`
- `sample_restrictions`: extract the exact condition including any mathematical formulation
and expected sample size.
BAD: "Among invalidated respondents"
GOOD: "Among invalidated respondents, defined as  $\text{sgn}(g_i) \neq \text{sgn}(\gamma_{\hat{i}})$  (Section 4.1). N = 1,365."

## 5. Cross-reference resolution

When anything references content elsewhere ("see Appendix H", "controls listed in Table A3",
"as defined in Section 2"), go to that location and extract the actual content. Never leave
unresolved references--- the replicator cannot see the paper. This is especially important for
control variable lists: find the appendix and enumerate every variable individually.

## 6. Variable completeness

Extract ALL variables: dependent, independent, controls, instruments, weights, and constructed
variables. For each, combine information from all locations in the paper (methodology section,
table notes, appendices). The goal: a replicator should be able to construct every variable
exactly as the authors did based solely on your extraction.

## 7. Table and figure precision

- Copy EXACT column headers and row labels as printed. Count columns and rows precisely.
- Include panel structure (Panel A / Panel B) when present.
- Copy full captions and table notes (excluding notes about specific coefficient values).
- For figures: note approximate axis ranges, reference lines, and line style conventions.
- If different tables use different datasets, specify the data source per table.""

EXTRACTION_USER_PROMPT = """Extract the methodology from this paper. Follow the system instructions.

## Detected Table Captions:
{table_captions}

## Detected Figure Captions:
{figure_captions}

## Paper Text:
{paper_text}"""

TEMPLATE_GENERATION_SYSTEM_PROMPT = """You are a structural template specialist for academic paper tables and figures.

Generate templates that faithfully reproduce the exact layout of tables and figures from a paper.

## Table templates

- EXACTLY the same number of columns and rows as the original table.
- Use the EXACT column headers and row labels from the paper.
- Cell content rules (check the ORIGINAL TABLE in the paper for each cell):
  - **XXX**: computed result (coefficient, statistic, count, mean, etc.)
  - *(XXX)*: standard error in parentheses
  - **Empty**: leave blank if blank in the original (do NOT use --- or placeholders)
  - **Literal text**: copy "Yes", "No", checkmarks, or labels as-is
- Include panel headers (Panel A, Panel B) as spanning rows when present.
- Include the full caption above the table.

## Figure templates

- Matplotlib code skeleton that makes the intended visual style unambiguous.
- Use the correct plot function: `ax.plot()`, `ax.bar()`, `ax.hist()`, `ax.scatter()`, etc.
- Set line styles, markers, colors for each series (use `tab:blue`, `tab:orange`, etc.).
- Set approximate axis ranges from the paper (`ax.set_xlim()`, `ax.set_ylim()`).
- Add reference lines (axvline, axhline) if the paper has them.
- Use actual series names / legend labels from the paper.
- Include the full caption as the title.
- Set up subplots correctly if the figure has panels.
- NO actual data arrays--- use comments like `# TODO: fill with data`.
- The skeleton should be runnable (producing an empty styled plot) even without data.""

TEMPLATE_GENERATION_USER_PROMPT = """Generate structural templates for each table and figure.

```

Respond with valid JSON.

Go back to the ORIGINAL TABLE in the paper text and reproduce its structure EXACTLY.

Example regression table template (note how Variable C only appears in columns (3)-(4)):

```
| | (1) | (2) | (3) | (4) |
|---|---|---|---|---|
| Variable A | XXX | XXX | XXX | XXX |
| | (XXX) | (XXX) | (XXX) | (XXX) |
| Variable B | XXX | XXX | XXX | XXX |
| | (XXX) | (XXX) | (XXX) | (XXX) |
| Variable C | | | XXX | XXX |
| | | (XXX) | (XXX) |
| Controls | No | Yes | No | Yes |
| Observations | XXX | XXX | XXX | XXX |
```

Example figure skeleton:

```
```python
import matplotlib.pyplot as plt
fig, ax = plt.subplots(1, 1, figsize=(10, 6))
x_obj, density_obj = ... # TODO: compute KDE from data
ax.plot(x_obj, density_obj, color="tab:blue", linestyle="--", linewidth=2, label="Objective")
ax.set_xlabel("Net gain (EUR/year)", fontsize=12)
ax.set_ylabel("Density", fontsize=12)
ax.set_xlim(-1500, 1500)
ax.axvline(x=0, color="black", linestyle=":", linewidth=0.8)
ax.legend()
plt.tight_layout()
plt.savefig("figure_X.png", dpi=150)
```

## Paper Text (excerpt):
{paper_text}

## Extracted Tables:
{table_specs_json}

## Extracted Figures:
{figure_specs_json}"""

# -----
# Non-replicable figure types (code-level safety net)
# -----

NON_REPLICABLE_KEYWORDS = [
    "flow diagram", "flow chart", "flowchart", "conceptual framework",
    "conceptual diagram", "diagram", "schematic", "screenshot", "photo",
    "photograph", "timeline",
]
]
```

C2. Reproduction Prompt

Replication Task

You are given a methodological summary of a research paper and its associated dataset. Your goal is to replicate the paper's empirical results using only the methodology description below and the data. You do not have access to the original paper, its results, or any replication code.

Data

The dataset is located at: `{data_filename}`

Methodological Summary

Paper: {title} (ID: {paper_id})

Research Questions:
{research_questions}

```

**Data Description**: {data_description}

**Data Context**: {data_context}

{data_source}{time_period}

**Data Processing Steps**:
{processing_steps}

{items_section}

## Constraints

You are in an isolated workspace for fair benchmarking.

1. **Workspace only.** Only read and write files inside this workspace directory.
   Do not access files outside of it or navigate to parent directories.
   **Exception:** the `data/` directory may be a symlink pointing outside the
   workspace -- this is intentional and you should use it freely as your dataset.

2. **No searching for the paper.** Do not search the internet for this paper,
   its authors, its published results, or any replication code or packages.
   Do not search for prior replication attempts.

3. **No searching for results.** Do not look up expected coefficients, effect
   sizes, tables, or figures from this paper. Your replication must be derived
   entirely from the methodology summary and the data provided.

4. **Allowed web use.** You may search for Python library documentation
   (e.g. statsmodels, pandas, matplotlib) and general statistical methods.

5. **Work independently.** Base your replication only on the methodology
   description in this file and the dataset.

## Instructions

1. **Quick data check**: Inspect the data files to confirm column names and
   basic structure. The methodology summary above describes the variables in
   detail, so a brief check should suffice.

2. **Write `prepare_data.py`**: Load and clean the data following the processing
   steps described above. All scripts can import from this module.

3. **Write and execute one script at a time**: For each item:
   a. Write the script (see output filename specified for each item above)
   b. Execute it
   c. Fix any errors immediately
   d. Move on to the next item once the output file is verified

{table_instructions}{figure_instructions}

4. **R packages**: If the paper's methodology requires R-specific packages
   (e.g. for specialized estimators), you can call R from Python using `rpy2`.

5. Execute every script and verify the output file exists.

6. Save all outputs in the current working directory.

**Reasonable assumptions.** Where the methodology description is incomplete or
   ambiguous, you are free to make reasonable assumptions based on common practice
   in the field. Document your assumptions briefly in comments.

Focus on substance and accuracy. Match the described methodology as closely as
   possible, including sample restrictions, variable transformations, and
   statistical specifications.

```

Note that, while Prompt C.2 specifies that web use is allowed for researching “Python library documentation” and “general statistical methods,” web search was in fact disabled programmatically at the scaffold level.

C3. Explanation - Extraction of Divergences Prompt

```
You are an expert at tracing numerical replication failures back to their root code discrepancies.

Your working directory contains two subfolders:
original_code/ -- original {original_language} replication package (code + data directory)
agent_code/    -- agent-generated Python replicator

## Target failure
Item: {item_id} | Grade: {grade} | {n_failed} of {n_total} cells failed

{already_attributed_block}Failed cells still to explain (sorted by [% difference]):
{cell_table}

## Your tasks

1. DATA AVAILABILITY --- Check whether the data files required to produce {item_id} exist in
the data directory inside original_code/. If any required input file is absent from the
replication package, set data_available = "missing" and describe the missing file(s).
Otherwise set data_available = "available".

2. AGENT CODE --- In agent_code/*.py, find the code that computes values for the REMAINING
cells listed above (ignore already-attributed sections).
If no such code exists, set agent_behavior = "<description of what is missing>".

3. ORIGINAL CODE --- In original_code/{original_file_glob}, find the equivalent {original_language} code.

4. DISCREPANCIES --- For the REMAINING cells only, identify all distinct root causes.
Most outputs have one, but multi-panel tables with different specifications may have
several (e.g. columns 1--3 OLS wrong clustering; columns 4--6 IV wrong instrument).
If all remaining cells are already explained, return an empty divergences array.
List one entry per distinct root cause --- do NOT split trivial variations.

5. SECTION MAPPING --- For each discrepancy, list which rows/columns/panels it explains
(e.g. "All columns", "Columns 1--3 (OLS)", "Panel B"). Maps many cells to one cause.
Also list each specific affected cell from the cell table above as
{"item_id": "{item_id}", "row_label": "...", "column_label": "..."} using the exact
labels shown. Include all cells you believe this discrepancy explains.

6. ALSO EXPLAINS --- For each discrepancy, list which other failures it also explains.
Each entry is either:
- A plain string "<item_id>" if the discrepancy explains the entire other item, OR
- An object {"item_id": "<item_id>", "sections": "<which part>"} if only partial.
Be conservative: only include items you are confident share the exact same root cause.

## Other unresolved failures
{remaining_block}

## Divergence taxonomy (use exactly one code)
S1 Wrong model specification -- incorrect FE, clustering level, or SE type
S2 Wrong estimator / inference -- wrong estimator (OLS vs IV) or missing inference step
S3 Data source substitution -- proxy used / required dataset absent from package
S4 Wrong sample restriction -- filter missing, wrong condition, or wrong subset
S5 Wrong variable construction -- outcome/predictor coded differently than reference
S6 Missing analysis component -- required step entirely omitted
S8 Wrong merge / transform logic -- wrong join type, key, duplicate handling, or reshape
S9 Wrong sequencing -- steps in wrong order, changing results
S0 Other -- does not fit any category above

## Severity
minor -- unlikely to materially affect point estimates or conclusions
medium -- could shift estimates noticeably; sign/significance probably stable
critical -- likely changes sign, significance, or core conclusion of a main result

## Output
Write a JSON object with a "divergences" array to {output_path}.
One entry per distinct root cause found in {item_id} (usually 1, occasionally 2--3).
Write ONLY valid JSON --- no markdown fences, no preamble, no comments.
Escape all special characters properly inside string values.

{{
  "divergences": [
    {{
```

```

    "id": {divergence_id}, (use sequential integers starting here for additional entries)
    "output": "{item_id}",
    "description": "<one sentence: what is different>",
    "original_behavior": "<what the original {original_language} code does on this specific point>",
    "agent_behavior": "<what Python does, or the ABSENT description>",
    "original_proof": "<short verbatim {original_proof_label}>",
    "agent_proof": "<short verbatim Python snippet, or ABSENT>",
    "original_location": {"file": "<{original_file_ext}>", "line": "<line or range>"},
    "agent_location": {"file": "<filename.py or ABSENT>", "line": "<line or ABSENT>"},
    "divergence_type": "<S0|S1|S2|S3|S4|S5|S6|S8|S9>",
    "severity": "<minor|medium|critical>",
    "data_available": "<available|missing>",
    "data_available_note": "<one sentence naming specific file(s) checked or found absent>",
    "explains_sections": ["<e.g. 'All columns', 'Columns 1-3 (OLS)', 'Panel B'>"],
    "affected_cells": [
    {"item_id": "{item_id}", "row_label": "<exact row label from the cell table>", "column_label": "<exact column label>"},
    ],
    "also_explains": ["<item_id>", {"item_id": "<item_id>", "sections": "<which part>"}]
  }
}
]]

```

C4.1. Explanation - Error source Original Paper vs Code

You are verifying whether the paper supports the original {original_language} code's behavior for a set of replication divergences.

Your working directory contains:

```

paper.pdf          --- the published paper
original_code_files/ --- the original {original_language} replication code (for reference only)

```

For each divergence you are given:

- `original_behavior`: what the {original_language} code does for this analysis step
- `original_proof`: the exact {original_language} code snippet implementing this behavior
- `original_location`: the file and line number in the original code

The original code is already provided --- you do NOT need to re-read the code files. Read paper.pdf to determine whether it explicitly states, implies, or is silent about the behavior described in `original_behavior`. Classify using exactly one of these four verdicts:

```

consistent = paper and original code explicitly agree on this specific point
contradicts = direct contradiction: the paper explicitly says X AND the original code explicitly does Y, where X ≠ Y
omission    = the original code implements X, but the paper does not mention X at all --- no contradiction, the paper simply doesn't cover this detail

```

Rules:

- "contradicts" requires both documents to explicitly address the same point with different answers. The paper being silent on an implementation detail is NEVER sufficient for "contradicts" --- use "omission" instead.
- "omission" means the original code (the upstream document here) specifies something the paper (the downstream document) doesn't address. Example: code computes monthly standardized values and averages them; paper describes the monthly formula but says nothing about aggregation → omission.
- "consistent" if both agree; prefer "omission" when one side is explicit and the other is silent.
- Your note must cite specific evidence: variable name, section heading, line number, or direct quote.
- You MUST produce exactly one of these three verdicts for every divergence id. Do not abstain or invent new labels.

EFFICIENCY: Do NOT read the entire paper exhaustively. Use the divergence description to identify the relevant section, check that section, form a well-founded hypothesis, cite the evidence, and move on.

C4.2. Explanation - Error source Original Paper vs Extractor

You are verifying whether the methodology summary accurately represents what the paper says, for a set of replication divergences.

Your working directory contains:

```

paper.pdf          --- the published paper

```

```
methodology_summary.json --- the summary passed to the replicator agent
```

For each divergence you are given:

- ``original_behavior``: what the original code does for this analysis step
- ``original_proof``: the exact code snippet (for context on what to look for)
- ``original_location``: the file and line number in the original code

The original code is already provided --- you do NOT need to read the code files. Read `paper.pdf` and `methodology_summary.json` to determine whether the paper states this behavior and whether the summary faithfully represents it. Classify using exactly one of these four verdicts:

- `consistent` = summary faithfully represents the paper on this point, OR both are silent (summary correctly captured the paper's silence)
- `contradicts` = direct contradiction: paper explicitly says X AND summary explicitly says $Y \neq X$
- `omission` = paper explicitly says X, but the summary does not mention X at all --- the summary dropped information the paper provided

Rules:

- "omission" means the paper (upstream) says something the summary (downstream) dropped. This is the most common failure mode for this check.
- "contradicts" requires both documents to explicitly address the same point differently --- not just one being silent where the other speaks.
- If neither the paper nor the summary addresses the point, mark "consistent" (they agree by silence --- the break is not here).
- The summary is NOT expected to be more detailed than the paper. If the paper states X at a conceptual level (e.g. "cluster at neighborhood level") and the summary conveys the same concept (e.g. "use neighborhood cluster identifiers"), mark "consistent" --- even if the summary does not reproduce the exact variable names or implementation details that appear in the original code but not in the paper. The standard is semantic faithfulness to what the paper says, not completeness of implementation detail.
- To distinguish "consistent" from "omission" at the conceptual level, ask: "Would a careful agent reading only the summary know that they should implement X (the concept the paper describes)?" If yes → "consistent" (the concept was conveyed, even if less specific). If no, because the summary says nothing about X at all → "omission". Note: "omission" does not require that the agent would actively do the wrong thing; it is sufficient that the summary dropped a concept the paper stated, leaving the agent without guidance on that point.
- Your note must cite specific evidence from both documents.
- You MUST produce exactly one of these three verdicts for every divergence id. Do not abstain or invent new labels.

EFFICIENCY: Do NOT read documents exhaustively. Use the divergence description to locate the relevant section in each document, check it, form a well-founded hypothesis, cite the evidence, and move on.

C4.3. Explanation - Error source Extractor vs Agent

You are verifying whether the agent's Python code implements what the methodology summary instructs, for a set of replication divergences.

Your working directory contains:

```
methodology_summary.json --- the summary passed to the replicator agent
agent_code/                --- the agent's Python replication code (for reference only)
```

For each divergence you are given:

- ``agent_behavior``: what the Python code actually does for this analysis step
- ``agent_proof``: the exact Python code snippet implementing this behavior
- ``agent_location``: the file and line number in `agent_code/`

The Python code is already provided --- you do NOT need to re-read `agent_code/`. Read `methodology_summary.json` to determine whether it explicitly instructs, implies, or is silent about the behavior described in ``agent_behavior``. Classify using exactly one of these four verdicts:

- `consistent` = agent follows the summary on this point, OR both are silent (agent correctly followed the summary's omission)
- `contradicts` = direct contradiction: summary explicitly says X AND agent code explicitly does $Y \neq X$
- `omission` = summary explicitly instructs X, but the agent does not implement X at all --- agent omitted something the summary described

Rules:

- Each divergence also includes a `description` field explaining what the step is about. Use this to find the relevant section in methodology_summary.json when agent_behavior is "ABSENT" or otherwise unclear.
- "omission" covers agent_behavior = "ABSENT": use `description` to identify what analysis step was supposed to be implemented, look it up in the summary, and if the summary describes it → omission. Only mark "consistent" if the summary is also silent on that step.
- "contradicts" requires the summary to explicitly say X AND the agent to do something explicitly different --- not just the agent being silent.
- IMPORTANT: if the summary does NOT mention this specific detail at all (e.g. the summary is silent about weighting, file choice, variable variant, clustering method, or any other implementation detail), mark "consistent" --- the agent cannot be expected to follow guidance that was not provided. The agent making its own reasonable choice on an unspecified detail is NOT a contradiction or omission.
- Your note must cite specific evidence from both the summary and the Python code.
- You MUST produce exactly one of these three verdicts for every divergence id. Do not abstain or invent new labels.

EFFICIENCY: Work from the provided code snippets and the summary text. Do NOT exhaustively search through all files. Form a well-founded hypothesis based on the evidence given, cite it, and move on.

C4.4. Explanation - Error source Data Missing

You are checking whether the data required by the original {original_language} code is available to a Python replicator, for a set of divergences.

Your working directory is the replication data directory. The original {original_code_description} are in {code_files_path}.

```
## Step 1 --- identify what the original code loads
```

For each divergence, read the original code files and find the data file(s) loaded for that analysis step.

```
## Step 2 --- distinguish raw inputs from code-constructed intermediates
```

Before checking whether a file exists, determine whether it is:

- (a) A RAW SOURCE FILE --- brought in from outside (survey data, official statistics, downloaded datasets). These must exist as-is for any replicator.
- (b) A CODE-CONSTRUCTED INTERMEDIATE --- a file that is itself created by the original code in the same replication package. A Python replicator would construct this from its own upstream steps, not load a pre-built file.

You can identify constructed intermediates by searching the original code for save/export commands that write that filename.

```
## Step 3 --- check availability
```

- If the required file is a RAW SOURCE FILE: check whether it exists in the current directory. If yes → available; if no → missing.
- If the required file is a CODE-CONSTRUCTED INTERMEDIATE: instead check whether the raw source file(s) that feed into its construction are present. If those raw sources exist → available (the agent can construct the intermediate from them); if the raw sources are also absent → missing.

available = the required data (raw or constructable from present sources) was accessible; the agent had what it needed and chose incorrectly
missing = the required raw data is absent; the agent could not have implemented this correctly regardless of effort

Every divergence must receive `available` or `missing` --- no other values.

In your note, name the file(s) you checked and state whether each is a raw source or a constructed intermediate.

```
\end{lstlisting}
```