

# Bridging Research Gaps With AI: Expertise and Disciplinary Mobility

Ghina Abdul Baki, Juan P. Aparicio, Bruno Barbarioli, Abel Brodeur,  
Lenka Fiala, Derek Mikola, David Valenta

2025-10-15

## 1 Abstract

We test whether providing large-language-model (LLM) assistance compresses performance gaps in research skills along two prespecified dimensions. Vertically, we study expertise tiers using a three-category grouping (undergraduate, graduate student, professor/researcher) that governs both randomization and analysis while retaining self-reported titles for descriptive context. Horizontally, we study cross-discipline performance in three quantitative social sciences (Economics, Political Science, Psychology) to assess whether AI reduces the penalty from working outside one’s primary discipline (out-of-discipline, OOD). Participants are randomized 1:1 to AI access (e.g. ChatGPT) versus human-only controls within the applicable expertise strata; tasks are assigned a single discipline tag, with undergraduates always receiving inside-discipline assignments and approximately the remaining participants randomly assigned to an outside-discipline paper. Outcome measures are organized into two families: coding skills (computational reproducibility and the detection of major/minor coding errors) and non-coding skills (referee appropriateness/overall scores and robustness-check execution). The study is designed to reveal whether AI acts as an equalizer within expertise tiers and whether it enables researchers to operate effectively across disciplines.

## 2 Registration and Funding

This pre-analysis plan is preregistered on the Open Science Framework (OSF) at <https://osf.io/dkfzt/>. Funded by the Social Sciences and Humanities Research Council. The locked PAP together with the analysis scripts used to generate the mock tables and figures are mirrored there; live study code and data will be added only after the analysis lock.

## 3 Background and Rationale

Prior multi-site “AI Replication Games” documented measurable AI effects on success and speed in reproduction tasks, while also revealing substantial variation across participants and events. Building on that foundation, we preregister two complementary dimensions of distributional impacts. First, the vertical dimension (expertise ladder): whether AI narrows performance gaps by disproportionately lifting less-experienced participants (equalizer) or widens gaps by enabling experts to better leverage tools (amplifier). Second, the horizontal dimension (cross-discipline): whether AI reduces the penalty from working outside one’s primary discipline (out-of-discipline, OOD) when

reproducing studies across Economics, Political Science, and Psychology. These questions matter for pedagogy, workforce development, and equity: vertical compression would support broad inclusion strategies; horizontal compression would support cross-field mobility and knowledge diffusion; in both cases, amplification would call for targeted training and governance to avoid widening disparities.

We keep tasks, instructions, and grading rubrics closely aligned with prior exercises to ensure comparability while tailoring the design to individual-level randomization within strata and a discipline-tagged task pool. The design isolates intent-to-treat effects within expertise tiers and introduces an orthogonal OOD contrast by allocating non-undergraduate participants to an outside-discipline article while keeping all undergraduates inside their primary field. This enables transparent tests of heterogeneous effects across tiers (vertical) and along the OOD dimension (horizontal). To support interpretability, we prespecify a compact outcome set split between coding skills (computational reproducibility; identification of major and minor coding errors) and non-coding skills (referee appropriateness/overall judgments and robustness-check execution) alongside a small list of precision-enhancing controls; discipline fixed effects are nested within article fixed effects and are thus absorbed in all models.

Several design features anchor the study and deserve to be unpacked carefully. We begin with randomization: participants enter one of three expertise strata—undergraduate, graduate (master’s and PhD students together), or professor/researcher (postdocs and faculty). Within each stratum, participants are randomized 1:1 to either AI assistance or a human-only control condition. This stratification keeps comparisons within peer groups while preserving overall balance; it also fixes the reference point for the “vertical” estimands that trace how treatment effects evolve along the expertise ladder.

Task assignment introduces the horizontal dimension. Articles arrive with a single discipline tag (Economics, Political Science, or Psychology), enabling us to contrast performance inside versus outside a participant’s primary field. We explicitly model two regressions—vertical ( $\text{AI} \times \text{tier}$ ) and horizontal ( $\text{AI} \times \text{out-of-discipline}$ )—and we test each dimension on its own set of outcomes. This separation keeps the interpretation of compression tests transparent and avoids conflating discipline moves with tier moves.

Outcome measurement rounds out the design. Coding outcomes combine a binary indicator for successful computational reproducibility with the counts of major and minor coding errors flagged during replication. Non-coding outcomes focus on communication and judgment: human referees—blinded to treatment and both participants and paper discipline—and a parallel AI grader each record whether the submission is appropriate and assign an overall 0–5 score, while participants’ robustness logs record whether they execute at least one or at least two qualifying checks. Section “Grading Rubric” details the anchors that guide human judgments, while Appendix “AI Referee Prompt” reproduces the corresponding instructions for the AI grader, ensuring that communication quality is assessed consistently across arms. We also track participants’ robustness checks proposals and implementations as part of the non-coding skill set to contextualize main effects.

## 4 Research Questions and Hypotheses

This study addresses three linked questions. First, we ask whether AI assistance compresses or widens the expertise gradient by comparing treatment effects across the three pre-defined tiers; this provides the “vertical” lens that motivates the stratified randomization. Second, we investigate

whether AI attenuates the penalty from tackling a problem outside one’s primary discipline, thereby capturing “horizontal” mobility across Economics, Political Science, and Psychology. Third, we consider the intent-to-treat contrast averaged across all participants, which anchors the benchmark effect size against which heterogeneous responses are interpreted.

The corresponding hypotheses follow naturally. H1 posits that access to AI increases research and coding skills—raising the probability of a successful reproduction while lowering major and minor coding errors. H2 mirrors this logic across disciplines, anticipating that AI shrinks the outside-of-discipline penalty for all outcome variables. H3 posits that AI assistance improves coding and research performance.

Throughout the document we rely on several recurring definitions. “Reproduction success” means the participant’s final output matches the prespecified focal result. “Error detection” counts correctly identified coding errors and distinguishes between major issues—those that would alter the substantive result—and minor issues that affect presentation or reproducibility without changing the significance and magnitude of the estimate. “Referee report quality” focuses on the Appropriateness indicator (binary) and the overall 0–5 score derived from the referee rubric; the underlying rubric is available to judges but only these two statistics enter the preregistered estimands. “Robustness execution” records whether participants implement at least one or at least two robustness checks that meet the prespecified criteria. Finally, “out-of-discipline (OOD)” labels any case in which the task’s discipline tag differs from the participant’s self-reported primary discipline.

## 5 Experimental Design

We recruit participants across the three canonical tiers—undergraduates, graduate students (master’s and PhD), and professors/researchers (postdocs and faculty)—and observe each in a single, timed, one-day session. Randomization always occurs within these three strata while we continue to track individual titles for descriptive reporting. Every participant attempts to reproduce one prespecified result using their preferred software ecosystems (R, Stata, or Python) within a seven-hour working window. The window mirrors prior AI Replication Games, where most participants completed their submissions within seven hours while still allowing careful documentation. We randomly assign access to AI assistance (ChatGPT Plus with tools) within the relevant strata and events. Participants in the control arm pledge not to use AI tools. Human referees are blinded to both treatment and participants and paper discipline, and filenames/metadata that could reveal either are redacted.

Within each event, articles are randomly assigned from a pre-curated pool. Once participant rosters are finalized. After the random draw, the research team audits each assigned paper to confirm that (i) the replication package includes executable code together with a README or documentation, (ii) all requisite datasets and intermediate files are accessible without additional permissions, and (iii) the combined data footprint remains tractable for local execution (target: compressed package 500 MB and memory requirements within a standard 16 GB RAM laptop). Papers that fail any check are replaced before assignments are communicated to participants.

Task materials are version-controlled in `Papers/`, which stores a journal-level folder for each study together with the inventory workbook `Papers/papers.xlsx`. As of the lock, the workbook enumerates 15 studies (American Journal of Political Science = 5, American Economic Review: Applied Economics = 5, Psychological Science = 5) with DOI identifiers and replication-package URLs. Each journal folder contains the published article (`paper.pdf`) and the original replication package

supplied by the journal or data repository. This structure lets us share identical bundles across events while keeping provenance and updates transparent.

The design studies vertical compression across expertise and horizontal compression across disciplines in a unified framework (run simultaneously). Each article carries a single discipline tag—Economics, Political Science, or Psychology—and participants self-report a primary discipline at registration. The task pool spans all three fields at every event. Undergraduates always receive inside-discipline papers; among graduate students and professors/researchers (postdocs and faculty) we randomly select individuals to work outside their primary discipline (OOD), with the remainder kept inside. This prioritizes a fixed OOD exposure rate over per-cell balance. Event  $\times$  article fixed effects absorb site, tooling, and task heterogeneity; because each article maps to a single discipline, separate discipline indicators are redundant.

## 5.1 Pre-Event Training

Participants attend a 60-minute orientation the week before each event (slides in `Pre game/ai_research_webinar_codex_cli_v2.pdf`). The session focuses on logistics: how the games run, what deliverables to submit, the structure of the Excel tracking file, and expectations about documentation, timing, and referee reports. We provide a brief overview of ChatGPT Plus features so attendees understand the tools that treated participants will receive, but the emphasis is on workflow discipline, reproducibility reminders, and answering procedural questions. Randomization is revealed one day before each event so that ChatGPT Plus invitations can be activated for treated participants in advance.

## 5.2 Focus Groups

To contextualize quantitative findings, we will run four parallel focus-group sessions (six participants per group) immediately after the first event concludes. Groups are stratified by treatment status (AI-assisted vs. human-only) and, when numbers permit, by inside/outside-discipline assignments so that discussions surface arm-specific workflows and cross-discipline frictions. Each 60-minute session is moderated by Institute for Replication staff using the standardized guide in `focus_groups/focus_group_guide.md`; facilitators remind attendees of confidentiality expectations and collect recorded consent before beginning. Discussion notes and transcripts are coded with the companion qualitative codebook (`focus_groups/codebook.md`), which maps themes on motivations, preparation, AI use or workarounds, cross-discipline challenges, and suggestions for future waves. Insights from these sessions inform protocol refinements before scaling beyond the first event.

Design cells overview (what is crossed with what):

Table 1: Design factors and levels (crossed: Arm  $\times$  Tier  $\times$  OOD).

Factor	Levels	Notes
Arm	Human-Only; AI-Assisted	1:1 allocation within tier
Expertise tier	Undergraduate; Graduate; Professor/Researcher	Stratification (randomization blocks; Graduate = MA/PhD, Professor/Researcher = postdoc/faculty)
OOD status	Inside; Outside	Derived from participant vs task discipline

Randomization is implemented with a reproducible script and a fixed seed recorded in the registry; the assignment file is timestamped and stored read-only. Allocation is concealed until check-in, when the onsite coordinator reveals arm and task. No-shows remain in their assigned arm for intent-to-treat analyses. Replacements are permitted only before the event begins and are

re-randomized using the same stratum-specific seed. The same concealment and documentation protocol applies to inside/outside assignments. Any late swaps or deviations are logged prior to accessing outcomes.

We plan to enroll roughly 300 participants across multiple events. For power, we assume a plausible tier composition (more undergraduates than postdocs/professors) rather than equal counts by tier. Simulations suggest that, with baseline success gaps between undergraduates and professors of 15–20 percentage points and AI compressing roughly 40% of that gap, we achieve at least 80% power for the vertical interactions at  $\alpha = 0.05$ . For the horizontal dimension, we assume a baseline OOD penalty of 20–25 percentage points in the control arm and target detectable AI-induced reductions of about 8–10 percentage points, while preserving the same overall sample size. Throughout, we rely on heteroskedasticity-robust standard errors. Table 2 summarizes the core inputs; we will freeze any updates to these assumptions prior to registry lock.

The primary analysis set follows intent-to-treat principles and includes all randomized individuals with any outcome data. We do not impute outcomes. For covariates only, we handle missingness as follows: categorical covariates (tier, software, prior ChatGPT familiarity, and—where used descriptively—participant and task discipline) gain an explicit “Missing” category if needed; continuous covariates (years of coding) use within-stratum median imputation with a missingness indicator. We report outcome and covariate missingness by arm and verify robustness to listwise deletion.

Table 2: Prospective power and design inputs (pre-lock)

Quantity	Value
Participants (N)	300 (approximate)
Tier composition (assumed)	UG 35–40%, Graduate 35–40%, Professor/Researcher 20–25%
Discipline composition	≈50% Econ, 25% PolSci, 25% Psych
Allocation	1:1 within tier (AI vs Control)
Inside vs Outside discipline	Undergrads always inside
Variance estimator	heteroskedasticity-robust SEs
Vertical: control success (UG; P)	40% ; 55–60% (assumed)
Vertical: baseline gap (UG vs Prof)	15–20 pp (assumed)
Vertical: detectable compression	≈ 40% of gap @ 80% power, $\alpha = 0.05$
Vertical: baseline major/minor error rates	Major 0.40–0.50; Minor 0.80–1.00 (counts per participant, illustrative)
Horizontal: OOD penalty (control)	20–25 pp (assumed)
Horizontal: detectable AI reduction (Success)	≈ 8–10 pp @ 80% power, $\alpha = 0.05$
Horizontal: detectable AI reduction (Major errors)	≈ 0.15–0.20 fewer errors @ 80% power (illustrative)
Non-coding baseline (Appropriate; Overall 0–5)	45–55% ; mean 3.0–3.3 (assumed)
Robustness execution baseline	55–60% 1 check; 30–35% 2 checks (assumed)

## 6 Outcomes and Measurement

Participants receive a discipline-tagged article, a standardized instruction sheet, and a data/code package (as available). They are asked to reproduce a pre-specified focal result and to document their workflow. At the end of the seven-hour window, participants submit two items: (i) the standardized Excel workbook (`Reports/Replication_Log_Referee_Template.xlsx`) completed with all required tabs, and (ii) a short post-event SurveyMonkey questionnaire covering arm assignment, perceived impact of AI access, usage intensity (AI arm only), perceived disadvantage (human-only arm), and open-ended reflections on helpful/missing AI features. No additional files are required. Human graders (blinded to treatment and discipline) later use the Excel workbook and survey responses to score success, error detection, and referee outcomes using the pre-defined rubric. The

AI referee prompt in the Appendix mirrors the human rubric and is applied to the narrative fields captured in the workbook.

We standardize the workflow with the same Excel template. Sheet **00\_Main** captures session metadata (participant name, article identifiers, software, event, discipline tags, and the out-of-discipline flag). Sheet **01\_CodingErrors** provides the structured log for major/minor issues, including timestamp, affected element, narrative justification, evidence pointer, and a minor/major self-classification toggle (although it serves only as a reference). Sheet **02\_Computation** records whether the focal result was reproduced and when. Sheet **03\_Robustness** allocates parallel columns for up to two robustness checks—each with motives, specification changes, original versus reproduced estimates, and interpretation—enforcing the “maximum two” rule in the protocol. Sheet **04\_Referee\_Report** houses the scored rubric (appropriateness plus narrative fields that feed both human and AI assessments). The post-event SurveyMonkey form captures the additional perception questions listed above; responses are linked back to the workbook for grading and descriptive summaries. The instrument records name and email, asks participants to confirm their arm assignment, and then branches: AI-assisted participants rate usage intensity on a 0–100 slider, report perceived performance impact on a –100 to 100 slider, list helpful AI features, and indicate which AI tools were used (multi-select, with an open ‘Other’ option). Human-only participants rate perceived performance impact without AI on the same –100 to 100 slider and describe which AI features they would have found most helpful. All respondents can provide open-ended comments.

We organize outcomes into two families that mirror our hypotheses: coding skills and non-coding skills. Coding skills track whether participants reproduced the focal result and how effectively they detected major or minor coding errors. Non-coding skills cover referee-style communication (appropriateness and overall scores from both human and AI judges) together with robustness behavior (whether participants proposed and implemented at least one or two qualifying checks). These measures together characterize both the technical and communicative dimensions of performance.

We also record a small set of supplementary variables for context—such as the treatment-by-event-order interaction that captures learning dynamics—and retain self-reported years of coding and prior AI usage as moderators (entering as covariates in the main models). The table below consolidates definitions, types, and assessment sources for quick reference; all analyses and figures use these exact definitions.

Table 3: Outcomes and measurements (coding and non-coding skill families).

Family	Outcome	Type	Measurement / Assessment
Coding skill	Success	Binary	Core result reproduced by endline (yes/no).
	Error detection — major	Count	Number of major coding errors correctly identified (pre-defined rubric).
	Error detection — minor	Count	Number of minor coding errors correctly identified (pre-defined rubric).
	Referee appropriateness (human)	Binary	Appropriate vs. Not appropriate; human judges (blinded).
	Referee overall 0–5 (human)	0–5	Holistic 0–5 assessment; human score averaged across judges.
	Referee appropriateness (AI)	Binary	Appropriate vs. Not appropriate; AI judge using mirrored rubric.

	Referee overall 0–5 (AI)	0–5	Holistic 0–5 assessment from AI judge using mirrored rubric.
Non-coding skill	Robustness planned 1	Binary	Indicator for proposing at least one qualifying robustness check (rubric-based).
	Robustness planned 2	Binary	Indicator for proposing at least two qualifying robustness checks.
	Robustness implemented 1	Binary	Indicator for executing at least one qualifying robustness check with code and outputs.
	Robustness implemented 2	Binary	Indicator for executing at least two qualifying robustness checks with code and outputs.
Supplementary	Learning (treatment $\times$ event order)	Interaction	Assesses whether treatment effects evolve across repeated events.
Heterogeneity	Years of coding	Continuous	Self-reported; used for heterogeneity analysis (not an outcome).
	Prior AI usage	Categorical	Self-reported; included as precision control and used for heterogeneity analysis (not an outcome).

We pre-define classification of “major” versus “minor” errors and use a standardized grading rubric. As a preview of magnitudes, the table below (Table 4) reports simple means (and standard deviations) by arm, together with Welch tests for the difference in means. Figures then visualize the core outcomes that the analysis will focus on (Figure 1).

The figures referenced here are predefined mock-ups of the panels we will produce once data are available; they are rendered using synthetic or placeholder data solely to illustrate the visual layout and expected content. Their purpose is to communicate the intended presentation of our preregistered estimands rather than to reveal any substantive pattern at this stage.

Table 4: Comparison of Human-Only and AI-Assisted Metrics

Variable	Human-Only	AI-Assisted	Human-Only vs AI-Assisted
Reproduction	0.468 (0.502)	0.505 (0.502)	-0.036 [0.627]
Number of minor errors	0.367 (0.644)	0.286 (0.514)	0.081 [0.357]
Number of major errors	0.139 (0.348)	0.200 (0.508)	-0.061 [0.337]
Planned 1 robustness check	0.329 (0.473)	0.438 (0.499)	-0.109 [0.132]
Planned 2 robustness checks	0.165 (0.373)	0.162 (0.370)	0.003 [0.962]
Implemented 1 robustness check	0.215 (0.414)	0.324 (0.470)	-0.109 [0.098]
Implemented 2 robustness checks	0.089 (0.286)	0.105 (0.308)	-0.016 [0.714]
Appropriate (human)	0.367 (0.485)	0.429 (0.497)	-0.061 [0.401]
Overall 0–5 (human)	3.093 (0.850)	3.226 (0.787)	-0.133 [0.281]
Appropriate (AI)	0.481 (0.503)	0.448 (0.500)	0.033 [0.655]
Overall 0–5 (AI)	3.081 (0.947)	3.211 (0.831)	-0.130 [0.334]

*Note:* Columns 2–3 present means and standard deviations in parentheses for the two arms; column 4 presents the difference in means (Human-Only – AI-Assisted) and two-sided Welch p-values in brackets.

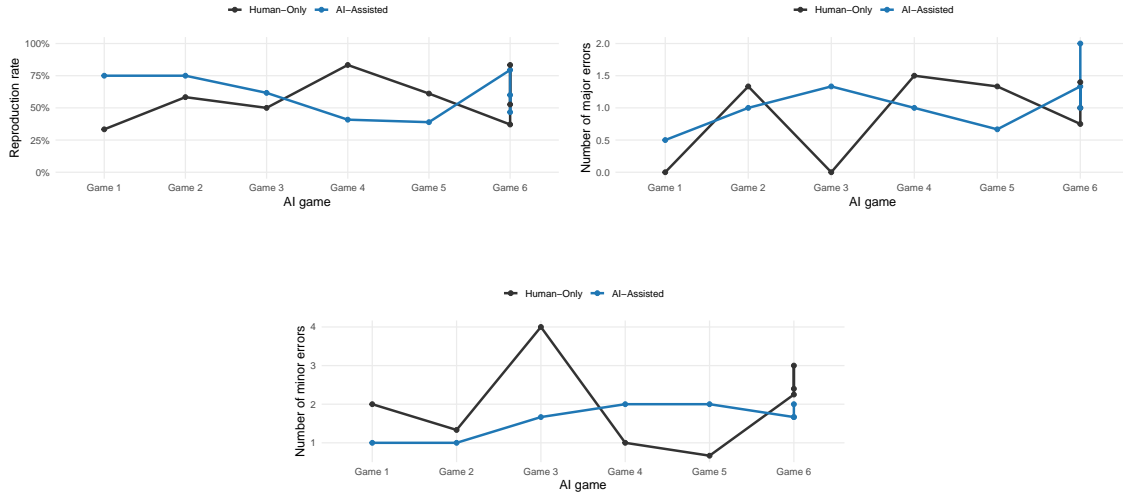


Figure 1: Coding skill outcomes (levels): reproduction and coding errors. Notes: Panels display reproduction rates (raw) alongside major and minor error counts.



## 6.1 Grading Rubric (Operational Definitions)

Coding skills - Success (binary): “Yes” when the participant’s final output matches the pre-specified focal result within documented tolerance (numeric threshold or visual equivalence for figures), with code that runs cleanly to produce the output. Partial or alternative results do not count. - Major error (count): Any coding/data/model issue that, once corrected, changes the focal result’s sign, statistical significance, or substantive interpretation; for example a wrong sample filter or an omitted transformation that alters the reported effect size. Issues that could matter but do not demonstrably affect the focal estimate are coded as minor. - Minor error (count): Formatting, non-substantive code, or reproducibility issues that do not change the focal result. Examples: mislabeled axes; non-deterministic seed; inefficient code without impact on estimate.

Non-coding skills - Appropriateness (binary): “Appropriate” when the report accurately identifies whether the reproduced evidence supports the original claim, flags substantive issues if present, and substantiates claims with code/snippets or references to outputs. - Overall (0–5): Holistic assessment of the referee report quality (not a mechanical average of component scores). Anchors follow the referee rubric’s 0–5 scale, where 0 denotes missing/incorrect and 5 denotes comprehensive, well-argued feedback. - Planned robustness checks: Indicators capturing whether the participant proposes at least one or at least two robustness checks that satisfy the prespecified rubric (clear motivation, documented specification change, comparison to the focal estimate, and concise interpretation). - Implemented robustness checks: Indicators capturing whether the participant successfully executes at least one or at least two qualifying checks, documented with code and resulting estimates.

Human graders apply these definitions blind to treatment and discipline; disagreements are reconciled via consensus. The same rubric is prompted to the AI judge for the AI-referee outcomes.

We will interpret Figure 1 as the primary visualization of coding outcomes (reproduction and error rates) with accompanying usage context. Difference-style and cumulative milestone panels are omitted to focus attention on the preregistered estimands.

## 7 Controls (Covariates and Stratification)

Table 5 summarizes the covariates and fixed effects we pre-specify for precision and design alignment. We use a compact, stable set that mirrors the randomization scheme and absorbs systematic heterogeneity without overfitting.

We include a small, pre-specified set of controls to improve precision and absorb systematic differences that are not of direct interest. Stratification by expertise tier (Undergraduate, Graduate, Professor/Researcher) reflects our design and is included as dummies in the main specification so that treatment effects are identified within tier; the interaction with treatment captures heterogeneous effects along the expertise ladder. Event and article fixed effects absorb site- and task-specific differences. Software indicators (R/Stata/Python) capture baseline workflow differences across toolchains. Finally, self-reported years of coding and prior AI familiarity improve precision and help stabilize estimates across events. Years of coding also serves as the pre-specified continuous moderator in secondary analyses; the resulting interaction terms do not alter the main ITT estimands.

Table 5: Controls, fixed effects, and moderators (prespecified).

Variable	Role	Type	Coding_or_Levels	Notes
Expertise tier	Stratification; control	Categorical	Undergrad, Graduate, Professor/Researcher	Tier dummies in main specs; Graduate = MA/PhD, Professor/Researcher = postdoc/faculty; interacted with AI (heterogeneity).
Event	Fixed effect	Categorical	One FE per event	Absorbs site/time differences (not a parameter of interest).
Article	Fixed effect	Categorical	One FE per task/article	Absorbs task-specific difficulty/fit (not a parameter).
Software	Control	Categorical	R, Stata, Python	Preferred software indicator (workflow baseline).
Years of coding	Control; moderator	Continuous	Self-reported years	Improves precision; interacted with AI in supplementary analyses (yrs×AI).
Prior ChatGPT familiarity	Control	Categorical	None, Some, Heavy	Self-reported familiarity with ChatGPT/AI tools.
Variance estimator	Estimation setting	—	heteroskedasticity-robust SEs	Used for all reported models.

## 8 Statistical Analysis Plan

This section lays out our estimands and modeling approach for the two preregistered dimensions of interest—vertical (expertise tier) and horizontal (out-of-discipline, OOD)—and clarifies how outcome types map to link functions and fixed-effect structure. We begin with intent-to-treat (ITT) specifications and then describe the interaction terms that capture compression:  $\text{AI} \times \text{tier}$  for vertical and  $\text{AI} \times \text{OOD}$  for horizontal.  $\text{Event} \times \text{article}$  fixed effects absorb site/tool and task heterogeneity; all models report heteroskedasticity-robust standard errors.

Let  $Y_i$  denote a pre-specified outcome (reproduction success; counts of minor or major coding errors; referee appropriateness; referee overall score; robustness-check indicators). We estimate ITT effects in a tier-interaction framework that allows the AI effect to vary with expertise. Formally,

$$Y_i = \beta_0 + \beta_1 A_i + \sum_s \gamma_s \mathbb{1}\{\text{Tier}_i = s\} + \sum_s \delta_s A_i \mathbb{1}\{\text{Tier}_i = s\} + X_i' \theta + \lambda_{e(i) \times s(i)} + \varepsilon_i,$$

where  $A_i$  is the treatment indicator,  $\lambda_{e(i) \times s(i)}$  are event-by-software fixed effects, and  $X_i$  collects the prespecified controls (years of coding and prior AI familiarity). Binary outcomes (reproduction, referee appropriateness, robustness indicators) are estimated with probit models; continuous outcomes (overall 0–5 scores) use OLS; and counts (minor/major errors) use a Poisson model:

$$\mathbb{E}[Y_i \mid \cdot] = \exp \left( \beta_0 + \beta_1 A_i + \sum_s \gamma_s \mathbb{1}\{\text{Tier}_i = s\} + \sum_s \delta_s A_i \mathbb{1}\{\text{Tier}_i = s\} + X_i' \theta + \lambda_{e(i) \times s(i)} \right).$$

Horizontal compression estimands. Let  $D_i$  denote an indicator for receiving a task outside the participant’s primary discipline (OOD). We test whether AI reduces the OOD penalty by augmenting

the main models with  $D_i$  and  $A_i \times D_i$ :  $Y_i = \dots + \beta_2 D_i + \beta_3 (A_i \times D_i) + \varepsilon_i$ . The coefficient  $\beta_3$  captures horizontal compression (a smaller penalty, or a gain, when outside). Article indicators absorb discipline effects because tasks carry a single discipline tag.

Interpretation and outcome-specific notes. In all tables, the “AI-Assisted” row reports the ITT contrast for participants in the baseline tier or inside-discipline case; the interaction rows report how that contrast varies along tiers (vertical) or OOD (horizontal). For counts, we check over-dispersion (deviance/df) and, if material, re-estimate with a negative binomial link as a robustness check.

To assess the success of randomization and support the modeling choices, we report balance on all individual-level controls used in the specification. The table below (Table 6) shows means (standard deviations) by arm and Welch tests for the difference in means.

Variable	Human-Only	AI-Assisted	Human-Only vs AI-Assisted
Years of coding	4.787 (2.775)	4.675 (2.819)	0.111 [0.747]
Tier: Undergraduate	0.218 (0.414)	0.237 (0.427)	-0.019 [0.720]
Tier: Graduate	0.429 (0.497)	0.435 (0.498)	-0.007 [0.915]
Tier: Professor/Researcher	0.353 (0.480)	0.328 (0.471)	0.025 [0.668]
Software: R	0.353 (0.480)	0.344 (0.477)	0.010 [0.867]
Software: Stata	0.368 (0.484)	0.344 (0.477)	0.025 [0.674]
Software: Python	0.278 (0.450)	0.313 (0.465)	-0.035 [0.538]
Prior ChatGPT familiarity: None	0.391 (0.490)	0.344 (0.477)	0.047 [0.426]
Prior ChatGPT familiarity: Some	0.323 (0.470)	0.359 (0.481)	-0.035 [0.545]
Prior ChatGPT familiarity: Heavy	0.286 (0.453)	0.298 (0.459)	-0.012 [0.831]

*Note:* Means and standard deviations in parentheses by arm; difference column shows Human-Only – AI-Assisted and two-sided Welch  $p$ -values in brackets. All variables are individual-level controls used in the models.

Supplementary analyses follow two paths. First, we replace tier dummies with a continuous moderator (years of coding) interacted with treatment to trace a dose–response. Second, within the AI arm only, we summarize SurveyMonkey responses on perceived AI usage intensity and performance impact; these descriptive measures help interpret treatment effects but do not alter the main ITT estimands. We also consider event-order interactions to gauge learning across events.

## 9 Results

For clarity and symmetry, we present vertical (tier-based) estimates first, followed by parallel horizontal (AI  $\times$  out-of-discipline) estimates using the same outcomes and table layout. Figures

and descriptives mirror this sequence.

We present results in two blocks that align with the preregistered outcome families—coding skills and non-coding skills—with vertical (tier) estimates first and horizontal (OOD) contrasts second. Throughout, standard errors rely on heteroskedasticity-robust corrections, and the coefficients are displayed with confidence intervals plus the pre-specified compression test for the interaction terms.

Table 7: Coding outcomes: vertical ITT effects with tier interactions. heteroskedasticity-robust standard errors.

	(1) Reproduction	(2) Minor errors	(3) Major errors
AI-Assisted	0.187 ( 0.119) [-0.047; 0.421]	-0.071 ( 0.465) [-0.982; 0.840]	1.132* ( 0.658) [-0.158; 2.422]
AI $\times$ Graduate	-0.170 ( 0.152) [-0.470; 0.130]	0.186 ( 0.590) [-0.971; 1.343]	-1.879* ( 1.003) [-3.844; 0.087]
AI $\times$ Professor/Researcher	0.032 ( 0.160) [-0.283; 0.346]	0.259 ( 0.569) [-0.856; 1.374]	-1.462* ( 0.810) [-3.049; 0.124]
Controls	✓	✓	✓
Mean of dep. var	0.485	0.356	0.159
p-val (Monotonic compression)	0.918	0.436	0.975
Obs.	264	264	264

*Note: Standard errors in parentheses; confidence intervals in brackets.*

Controls: Event $\times$ article FE; years of coding; software; prior AI familiarity.

Compression (monotonic): one-sided p-value for increasing effects across the three tiers (baseline: Undergraduate).

\* $p < 0.10$ , \*\* $p < 0.05$ , \*\*\* $p < 0.01$

Table 8: Coding outcomes: horizontal AI  $\times$  outside-of-discipline effects. heteroskedasticity-robust standard errors.

	(1) Reproduction	(2) Minor errors	(3) Major errors
AI-Assisted	0.339*** ( 0.096) [ 0.151; 0.528]	0.156 ( 0.309) [-0.450; 0.761]	0.230 ( 0.505) [-0.761; 1.221]
AI $\times$ Outside-discipline	-0.339*** ( 0.121) [-0.577; -0.100]	-0.154 ( 0.425) [-0.987; 0.680]	-0.382 ( 0.652) [-1.659; 0.895]
Controls	✓	✓	✓
Mean of dep. var	0.485	0.356	0.159
Obs.	264	264	264

*Note: Standard errors in parentheses; confidence intervals in brackets.*

Controls: Event $\times$ article FE; years of coding; software; prior AI familiarity.

\* $p < 0.10$ , \*\* $p < 0.05$ , \*\*\* $p < 0.01$

Table 9: Non-coding outcomes: vertical ITT effects on referee judgments and robustness execution. heteroskedasticity-robust standard errors.

	(1)	(2)	(3)	(4)	(5)	(6)	(7)	(8)
	Appropriate (human)	Overall 0-5 (human)	Appropriate (AI)	Overall 0-5 (AI)	At least 1 check	At least 2 checks	Implemented 1	Implemented 2
AI-Assisted	0.044 ( 0.129) [-0.210; 0.297]	0.210 ( 0.226) [-0.236; 0.655]	-0.102 ( 0.133) [-0.364; 0.159]	0.159 ( 0.243) [-0.321; 0.639]	-0.052 ( 0.131) [-0.311; 0.207]	0.176* ( 0.091) [-0.002; 0.354]	0.008 ( 0.114) [-0.216; 0.232]	0.149** ( 0.068) [ 0.015; 0.284]
AI × Graduate	0.059 ( 0.165) [-0.267; 0.384]	0.243 ( 0.264) [-0.277; 0.763]	0.145 ( 0.166) [-0.183; 0.473]	0.343 ( 0.280) [-0.210; 0.895]	0.118 ( 0.160) [-0.198; 0.434]	0.155 ( 0.117) [-0.386; 0.076]	0.131 ( 0.140) [-0.145; 0.406]	-0.063 ( 0.094) [-0.249; 0.122]
AI × Professor/Researcher	-0.145 ( 0.166) [-0.472; 0.181]	-0.079 ( 0.281) [-0.633; 0.476]	-0.043 ( 0.172) [-0.382; 0.295]	-0.026 ( 0.295) [-0.607; 0.555]	0.302* ( 0.170) [-0.032; 0.637]	-0.048 ( 0.122) [-0.288; 0.193]	0.209 ( 0.157) [-0.101; 0.518]	-0.133 ( 0.095) [-0.320; 0.053]
Controls	✓	✓	✓	✓	✓	✓	✓	✓
Mean of dep. var	0.443	3.178	0.462	3.190	0.360	0.170	0.254	0.110
p-val (Monotonic compression)	0.927	0.932	0.906	0.949	0.270	0.953	0.306	0.918
Obs.	264	264	264	264	264	264	264	264

Note: Standard errors in parentheses; confidence intervals in brackets.  
Controls: Event×article FE; years of coding; software; prior AI familiarity.  
Compression (monotonic): one-sided p-value for increasing effects across the three tiers (baseline: Undergraduate).  
\*p < 0.10, \*\*p < 0.05, \*\*\*p < 0.01

Table 10: Non-coding outcomes: horizontal AI × outside-of-discipline effects. heteroskedasticity-robust standard errors.

	(1)	(2)	(3)	(4)	(5)	(6)	(7)	(8)
	Appropriate (human)	Overall 0-5 (human)	Appropriate (AI)	Overall 0-5 (AI)	At least 1 check	At least 2 checks	Implemented 1	Implemented 2
AI-Assisted	0.019 ( 0.106) [-0.189; 0.228]	0.202 ( 0.164) [-0.121; 0.525]	-0.093 ( 0.101) [-0.293; 0.106]	0.258 ( 0.174) [-0.084; 0.600]	0.133 ( 0.105) [-0.074; 0.341]	0.204** ( 0.079) [ 0.048; 0.360]	0.227** ( 0.092) [ 0.045; 0.409]	0.119* ( 0.069) [-0.016; 0.255]
AI × Outside-discipline	-0.007 ( 0.130) [-0.263; 0.249]	0.127 ( 0.203) [-0.272; 0.526]	0.065 ( 0.129) [-0.190; 0.320]	0.058 ( 0.212) [-0.359; 0.475]	-0.047 ( 0.129) [-0.301; 0.208]	-0.173* ( 0.100) [-0.370; 0.023]	-0.139 ( 0.117) [-0.369; 0.091]	-0.070 ( 0.085) [-0.237; 0.097]
Controls	✓	✓	✓	✓	✓	✓	✓	✓
Mean of dep. var	0.443	3.178	0.462	3.190	0.360	0.170	0.254	0.110
Obs.	264	264	264	264	264	264	264	264

Note: Standard errors in parentheses; confidence intervals in brackets.  
Controls: Event×article FE; years of coding; software; prior AI familiarity.  
\*p < 0.10, \*\*p < 0.05, \*\*\*p < 0.01

## 10 Data Management, Documentation, and Ethics

We will publish de-identified participant-level data, code, and grading rubrics on OSF and GitHub upon acceptance (view-only earlier when necessary).

Static study materials live inside this repository and will be version-locked before launch: curated task bundles in **Papers/**, the pre-game training deck in **Pre game/**, and the reporting workbook in **Reports/**. Repository history provides provenance for any updates (for example, refreshed replication packages or revised slides).

Audio recordings and notes from the focus groups are stored on encrypted I4R drives with filenames keyed to anonymous participant IDs and session arm (AI vs. control). Transcription software with local processing is used where feasible; otherwise, trained staff transcribe manually. Only the qualitative analysis team accesses the raw audio. De-identified excerpts linked to codebook categories are released alongside the quantitative replication package after redacting personally identifiable information and any mention of other participants’ outputs.

## 11 Appendix

### 11.1 AI Referee Prompt

The AI grader receives a standalone instruction each time it evaluates a submission. The block below reproduces the full prompt, including the required output structure and scoring anchors.

**Role and objective** You are “AI Referee,” an expert evaluator for the AI Replication Games.

Your task is to read a participant’s workbook and produce referee-style assessments that mirror the human rubric.

**Context** - You are judging a replication of a published empirical result. - The participant had access to the original paper, its replication package, and a seven-hour work window. - The materials you receive include the participant’s referee report (sheet `04_Referee_Report`), their reproducibility notes and robustness log, and any supporting tables or figures they produced.

**General instructions** - Work independently—do not infer facts that are absent from the materials. - Keep the narrative within 1,000 words and organize it under the sections listed below; fold any code or data blockers into those sections instead of creating new headings. - Base every judgment on the participant’s reasoning and evidence; parroting the source paper without evaluation should be penalized. - Highlight robustness checks only when they satisfy the rubric (motivation, documented specification change, comparison to the focal estimate, interpretation); otherwise describe the gap. - Use the full 0–5 range for the overall score (no half points). - Provide two- to four-sentence justifications for each scored dimension, citing the relevant portions of the workbook (quote, paraphrase, or reference to a table or figure). - If critical evidence is missing, record that as a weakness, reflect it in the score, and flag it in the notes.

**Evaluation tasks** 1. Appropriateness (binary). Return “Appropriate” if the participant correctly characterizes whether the reproduced evidence supports the original claim, flags material discrepancies, and grounds the conclusion in the submitted code or outputs. Otherwise return “Not appropriate”. 2. Overall assessment (0–5). Provide a holistic rating that reflects the usefulness of the referee narrative to an editor, weighing accuracy, diagnostic insight, actionable guidance, and clarity. This is not a mechanical average.

**Anchors for the 0–5 scale (task 2)** - 0 = Missing, wholly incorrect, or incoherent. - 1 = Substantially inaccurate, with major misunderstandings that would mislead an editor. - 2 = Partially correct but incomplete or containing notable inaccuracies that reduce usefulness. - 3 = Adequate: correct core points with limited depth or weak justification. - 4 = Strong: accurate, well-supported, and actionable with only minor omissions. - 5 = Exceptional: precise, comprehensive, and offering original insight or particularly valuable guidance.

**Report structure** Produce prose for the following sections. Each section should be a short paragraph that references the participant’s evidence and names any missing components. - Summary — question, setting, methods, and focal result (including whether it was reproduced and any tolerance issues). - Design and Identification — randomization, balance, attrition, and identification threats. - Robustness — up to two checks with rationale, specification or sample changes, and how the results compare with the focal estimate; note explicitly if required elements are absent. - Ethical and Transparency — preregistration, data access, and consent or IRB references, if any. - Overall Assessment — recommendation to editors/authors synthesizing strengths and weaknesses.

**Output format** Return JSON with the following fields: `appropriateness` (“Appropriate” or “Not appropriate”), `overall`, `justifications`, `report_sections`, and `notes`. The `justifications` object must contain keys for `appropriateness` and `overall`. The `report_sections` object must contain keys `summary`, `design_identification`, `robustness`, `ethical_transparency`, and `overall_assessment`. Include a single string in `notes` describing any missing information, suspected hallucinations, or uncertainties.

```
{  
  "appropriateness": "Appropriate",
```

```

"overall": 4,
"justifications": {
  "appropriateness": "...",
  "overall": "..."
},
"report_sections": {
  "summary": "...",
  "design_identification": "...",
  "robustness": "...",
  "ethical_transparency": "...",
  "overall_assessment": "..."
},
"notes": "..."
}

```

**Quality control** - Double-check that every field is filled. - If evidence is missing for a dimension: - set the score to the literal value `null`, - explain the omission in the `notes` field, and - flag the gap in the corresponding justification or section. - Do not invent references or cite external material. - Use only the materials provided with the submission. - Return only the JSON object—no additional prose after the closing brace.

Table 11: Coding outcomes with years-of-coding moderator (interaction with AI).

	(1) Reproduction (yrs)	(2) Minor (yrs)	(3) Major (yrs)
AI-Assisted	0.045 ( 0.118) [-0.187; 0.277]	-0.908** ( 0.422) [-1.734; -0.082]	0.758 ( 0.755) [-0.722; 2.237]
AI × Years of coding	0.017 ( 0.022) [-0.027; 0.060]	0.207*** ( 0.071) [ 0.068; 0.346]	-0.153 ( 0.134) [-0.415; 0.109]
Controls	✓	✓	✓
Mean of dep. var	0.485	0.356	0.159
Obs.	264	264	264

*Note: Standard errors in parentheses; confidence intervals in brackets.*

Controls: Event×article FE; software; prior AI familiarity.

\* $p < 0.10$ , \*\* $p < 0.05$ , \*\*\* $p < 0.01$

Table 12: Non-coding outcomes with years-of-coding moderator (interaction with AI).

	(1) Appropriate (human) (yrs)	(2) Overall 0-5 (human) (yrs)	(3) Appropriate (AI) (yrs)	(4) Overall 0-5 (AI) (yrs)	(5) Implemented 1 (yrs)	(6) Implemented 2 (yrs)	(7) Planned 1 (yrs)	(8) Planned 2 (yrs)
AI-Assisted	-0.021 ( 0.119) [-0.255; 0.213]	0.685*** ( 0.196) [ 0.299; 1.070]	-0.078 ( 0.125) [-0.325; 0.169]	0.717*** ( 0.212) [ 0.300; 1.135]	0.017 ( 0.112) [-0.203; 0.238]	0.161** ( 0.065) [ 0.032; 0.290]	-0.089 ( 0.123) [-0.332; 0.153]	0.080 ( 0.082) [-0.081; 0.242]
AI × Years of coding	0.008 ( 0.022) [-0.035; 0.052]	-0.084** ( 0.034) [-0.152; -0.016]	0.005 ( 0.023) [-0.040; 0.051]	-0.088** ( 0.038) [-0.163; -0.013]	0.025 ( 0.021) [-0.017; 0.067]	-0.018 ( 0.013) [-0.043; 0.007]	0.040* ( 0.022) [-0.004; 0.084]	0.003 ( 0.016) [-0.030; 0.035]
Controls	✓	✓	✓	✓	✓	✓	✓	✓
Mean of dep. var	0.443	3.178	0.462	3.190	0.254	0.110	0.360	0.170
Obs.	264	264	264	264	264	264	264	264

*Note: Standard errors in parentheses; confidence intervals in brackets.*

Controls: Event×article FE; software; prior AI familiarity.

\* $p < 0.10$ , \*\* $p < 0.05$ , \*\*\* $p < 0.01$

## 12 References

References render from `references.bib`. We will cite prior AI Replication Games and related methodology upon registration finalization.