

# 1 Formalized scientific methodology enables rigorous 2 AI-conducted research across domains

Yanlin Zhang<sup>1,\*</sup> Jing ZHAO<sup>2</sup>

<sup>1</sup> Data Science and Analytics Thrust, Hong Kong University of Science and Technology (Guangzhou), Guangzhou,  
3 511453, China

<sup>2</sup> Shenzhen DataSpeak Intelligent Technology Co., Shenzhen 518000, Guangdong, China

\*Corresponding author. Email: [yanlinzhang@hkust-gz.edu.cn](mailto:yanlinzhang@hkust-gz.edu.cn)

## 4 Abstract

5 We formalize scientific methodology—the end-to-end process from question formulation to evidence-  
6 grounded writing—as a phase-gated research protocol with explicit return paths and persistent con-  
7 straints, and instantiate it for general-purpose language models as executable protocol specifications.  
8 The formalization decomposes methodology into three complementary layers: a procedural workflow, an  
9 integrity discipline, and project governance. Encoded as protocol specifications<sup>§</sup> and activated across the  
10 lifecycle, these constraints externalize planning and verification artifacts and make integrity-relevant in-  
11 terventions auditable. We validate the approach in six end-to-end projects, including a matched controlled  
12 study, where the same agent produced two complete papers with and without the protocol. Across do-  
13 mains, the protocol-constrained agent produced evidence-backed, auditable research outputs—including  
14 closed-form derivations, quantitative ablations that resolve modeling design choices, and algorithmic  
15 refactors that preserve the objective while changing the computational primitive. In population-genomic  
16 applications, it also recovered well-studied biological signals as validity checks, including known admix-  
17 ture targets in the 1000 Genomes Project and Neanderthal-introgressed immune loci on chromosome 21  
18 consistent with prior catalogs. In the controlled study, the protocol-free baseline could still produce  
19 a complete manuscript, but integrity-relevant risks were easier to introduce and harder to detect when  
20 constraints and artifacts were absent.

---

<sup>§</sup>Source code: <https://github.com/EvoClaw/amplify>.

## 21 Introduction

22 Doing good science requires two kinds of knowledge. The first—*what* is known in a field—is encoded in  
23 textbooks, papers, and databases. The second—*how* to generate reliable new knowledge—is the methodol-  
24 ogy that governs research practice: formulate research questions that matter, lock evaluation criteria before  
25 running experiments, report all results including failures, exclude alternative explanations before claiming  
26 mechanisms, ensure every assertion rests on specific evidence. This procedural and normative knowledge is  
27 what separates a publishable study from an exploratory exercise. Yet while much of it can in principle be  
28 articulated—and fragments have been codified in pre-registration protocols [1], registered reports [2], and  
29 reporting checklists [3, 4]—it is most commonly learned and refined through apprenticeship: PhD training,  
30 mentorship, lab culture, and the accumulated experience of reviewer feedback [5–7].

31 For agentic research, this apprenticeship-based transmission creates a more immediate problem: key  
32 methodological steps are often implicit, context-dependent, and enforced socially (through advisor feedback,  
33 lab norms, and peer review) rather than encoded as explicit, checkable rules. When these expectations are not  
34 externalized, an agent can complete a project while silently skipping steps that humans would normally insist  
35 on (e.g., freezing evaluation criteria before experiments, reporting negative results with equal prominence,  
36 or mapping claims to concrete evidence). The result is not an inability to produce manuscripts, but a gap in  
37 auditability and process discipline that makes failures harder to detect, comparisons harder to run fairly, and  
38 iteration harder to govern.

39 Scientific discovery follows a generative cycle: a researcher formulates a question, designs experiments  
40 to address it, collects observations, and then integrates those observations with existing knowledge to pro-  
41 duce new understanding. Large language models have now demonstrated expert-level knowledge across  
42 numerous scientific domains [8, 9]—they possess, in effect, the first ingredient of this cycle. What is less  
43 reliably expressed in autonomous settings is the second: the procedural knowledge of *how* to move from  
44 domain expertise to reliable new findings, including when to stop, verify, and backtrack. If LLMs can be  
45 taught to formulate questions, design and execute experiments, verify results, and reason from evidence  
46 to conclusions—much as a graduate student learns from an experienced mentor—then the combination of  
47 broad domain knowledge with sound methodology should enable genuine knowledge creation.

48 Artificial intelligence is already transforming science, but current AI systems encode scientific *content*  
49 and task *capability* rather than an end-to-end, auditable scientific *methodology*. At one end, domain-specific

50 systems achieve deep integration with particular fields: AlphaFold embeds protein physics [10, 11], au-  
51 tonomous laboratories encode reaction rules [12, 13], and biomedical platforms such as Biomni [14] orches-  
52 trate hundreds of specialized tools across 25 subfields to execute tasks from drug repurposing to molecular  
53 cloning. At the other, AI scientist systems pursue research autonomy at increasing scale: The AI Scien-  
54 tist [15] generates and peer-reviews machine learning papers, its successor produced the first AI-generated  
55 peer-reviewed publication [16], Kosmos [17] performs 12-hour autonomous discovery sessions equivalent to  
56 months of human research, Robin [18] achieved the first fully automated biological discovery, and Google’s  
57 AI co-scientist [19] uses multi-agent debate and tournament evolution to generate and rank research hy-  
58 potheses for biomedical discovery. Between these, bioinformatics agents [20, 21] automate standard compu-  
59 tational pipelines, while deep research products [22] synthesize literature at unprecedented scale. These sys-  
60 tems represent genuine and rapid progress in what AI can *do* in science. What often remains under-specified  
61 is *how* AI should do it: the persistent methodological constraints—evaluation immutability, complete re-  
62 porting, claim–evidence alignment, alternative-hypothesis exclusion—that distinguish reliable knowledge  
63 from plausible-sounding output. Without such constraints, capable systems can more easily accumulate  
64 integrity-relevant risks (e.g., metric drift, incomplete reporting, unverified references, and claim–evidence  
65 mismatches) [23, 24]. Existing experiment tracking tools [25] and FAIR data principles [26] address com-  
66 putational reproducibility, and multi-agent frameworks improve coordination [27, 28], but these tools do not  
67 encode an end-to-end research methodology as an executable set of norms that can be audited turn by turn.  
68 Moreover, as with human researchers, the most consequential—and most difficult—step in the entire cycle  
69 is formulating a good scientific question. A single LLM mining the literature often tends toward conser-  
70 vative, incremental questions; the methodology formalization we propose addresses this through structured  
71 multi-agent deliberation and explicit ideation strategies that generate divergent candidate directions before  
72 convergence.

73 Here we ask whether the articulable components of scientific methodology—the procedural rules, in-  
74 tegrity norms, and governance practices that experienced researchers follow—can be operationalized end-  
75 to-end as an executable, phase-gated protocol with persistent methodological constraints, and transferred to  
76 general-purpose AI agents as a modular specification they execute step by step. The key insight is twofold.  
77 First, the complete research process—from question formulation through direction validation, method de-  
78 sign, iterative experimentation, failure management, results integration, and evidence-grounded writing—  
79 can be decomposed into three complementary constraint layers (procedural workflow, integrity discipline,

80 and project governance) and specified as an explicit protocol rather than an ad hoc interaction. Second,  
81 across our end-to-end validation projects the protocol condition externalizes intermediate planning and veri-  
82 fication artifacts, enforces phase transitions and explicit backtracks when evidence is stale, and helps surface  
83 and correct claim–evidence mismatches relative to matched protocol-free runs. We applied this formal-  
84 ization to six end-to-end projects spanning diverse domains, without domain-specific modification: five  
85 protocol-constrained projects executed under the full methodology plus one controlled study that produced  
86 two complete papers under matched conditions with and without the protocol (Supplementary Papers 1–  
87 5 and 6A/6B). Across these projects, the protocol-constrained agent produced evidence-backed, auditable  
88 research outputs—including closed-form derivations, quantitative ablations that resolve modeling design  
89 choices, and algorithmic refactors that preserve the objective while changing the computational primitive—  
90 and recovered well-studied biological signals as validity checks. However, it also inherits well-known failure  
91 modes of large language models—including simple arithmetic or consistency errors that can be introduced  
92 during manuscript drafting and may only be surfaced through multi-round, adversarial review. These valida-  
93 tion runs suggest that the formalization itself—not merely the underlying model capability—materially con-  
94 tributes to methodological rigour and auditability, as matched comparison runs without constraints exhibited  
95 more frequent integrity-relevant risks. Human researchers may participate at validation gates to contribute  
96 scientific judgment that remains difficult to formalize (scientific taste, intuition about significance [5]), but  
97 the methodology functions as a self-consistent system: the interlocking constraints guide the AI through the  
98 full research lifecycle, with each layer catching distinct failure modes. We instantiate this formalization as  
99 an open-source tool, *Amplify*, which implements the protocol and produces auditable artifacts; all validation  
100 projects and the controlled study in this paper are executed via Amplify.

## 101 **Results**

102 We propose a formalized, executable research methodology for general-purpose AI agents, implemented  
103 as a phase-gated protocol with persistent constraints. The protocol structures the full lifecycle into seven  
104 phases with explicit return paths, maintains always-on integrity and governance invariants, and embeds  
105 role-specialized multi-agent deliberation checkpoints that can trigger revision and backtracking. In a tool-  
106 enabled environment (IDE with web search/browsing, literature reading, and a shell for execution), the  
107 agent is permitted to iteratively gather evidence, run analyses, verify claims against fresh computations

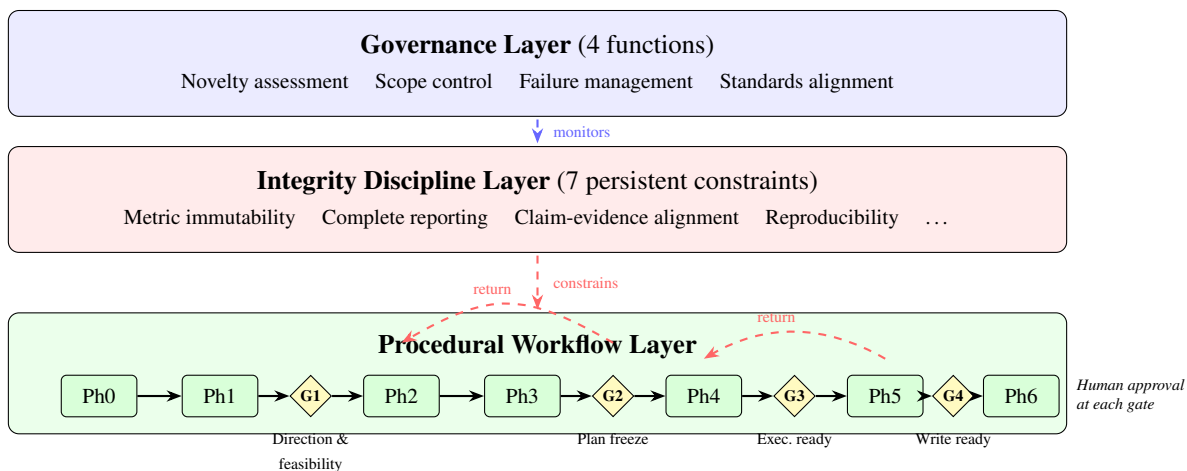


Figure 1: **Three-layer methodology formalization with validation gates.** Scientific methodology is decomposed into three complementary constraint layers: a procedural workflow (seven phase-gated phases with explicit return paths), an integrity discipline (seven persistent constraints active throughout), and a governance layer (four strategic oversight functions). Four mandatory gates (G1–G4, diamonds) require human approval at critical transitions, preserving scientific judgment while the formalization provides methodological discipline. The integrity and governance layers operate continuously across all phases.

108 and references, and revise artifacts until gate criteria are satisfied. We evaluate the approach through five  
109 end-to-end protocol-constrained projects spanning multiple domains and research types, and a matched con-  
110 trolled study (Project 6) in which the same task and environment produced two complete papers with and  
111 without the protocol; the six projects and manuscripts are summarized in Supplementary Table S1. Across  
112 the validation projects, the protocol externalizes intermediate planning and verification artifacts and makes  
113 integrity-relevant interventions auditable. In the controlled study, the protocol-free baseline could still pro-  
114 duce a complete manuscript, but the protocol condition more consistently enforced domain-target alignment  
115 and verification obligations by requiring those intermediate artifacts and return paths to be made explicit  
116 and checkable. Importantly, beyond process-level interventions, these projects produced evidence-backed  
117 research outputs that are auditable in the generated papers (Table 2), including analytical derivations, quan-  
118 titative ablations, algorithmic refactors, and recovery of well-studied biological signals as validity checks.

## 119 **Decomposing scientific methodology into functional layers**

120 To formalize scientific methodology, we first needed to determine whether its constituent practices share  
121 a structure amenable to decomposition. We analysed the methodological norms enforced by experienced  
122 researchers—as documented in integrity guidelines [29, 30], reporting standards [3, 4], and pre-registration

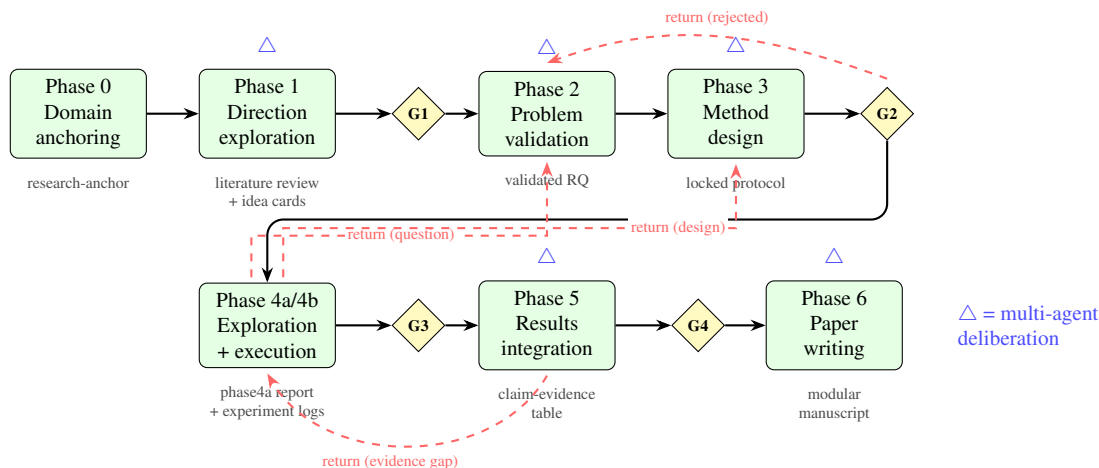


Figure 2: **Seven-phase procedural workflow with deliverables.** Each phase produces defined artifacts (grey text). Gates (diamonds) require human approval before progression. Blue triangles mark phases where multi-agent deliberation panels are deployed. Phase 4 is explicitly two-stage: Phase 4a exploratory probing that may trigger returns to refine the question or design, and Phase 4b full execution. Red dashed arrows show return paths triggered by governance interventions (e.g., plan rejection at G2), Phase 4a findings (question/design refinement), or evidence gaps identified during integration (Ph5→Ph4). The workflow enforces temporal ordering: evaluation protocols are locked at G2 before experiments begin, and evidence sufficiency is verified at G4 before writing begins.

123 frameworks [1, 2]—and observed that they cluster around three distinct functional concerns: *what to do*  
 124 *next* (procedural sequencing), *what constraints must always hold* (integrity norms), and *whether the project*  
 125 *should continue* (strategic governance).

126 We use this three-part organization as a practical decomposition: each concern can be specified and en-  
 127 forced via a corresponding constraint layer. This separation is not merely a convenient taxonomy but reflects  
 128 how methodology functions in practice. In well-functioning human research groups, procedural planning  
 129 (the PI’s research plan), integrity enforcement (lab culture, peer norms), and strategic oversight (committee  
 130 review, self-assessment) are typically carried out by different actors through different mechanisms [7]. We  
 131 formalized each concern as a separate constraint layer applicable to AI agents: a **procedural workflow** of  
 132 seven phase-gated phases with explicit return paths, an **integrity discipline** of seven persistent constraints,  
 133 and a **governance** layer of four strategic functions (Fig. 1).

134 We evaluate this decomposition through cross-domain application, and organize the Results below by  
 135 mapping failure modes observed when constraints are absent to the layer(s) designed to prevent them.

## 136 **Procedural workflow mitigates the risk of collapsing research into undifferentiated genera-** 137 **tion**

138 The procedural workflow encodes seven phase-gated phases (domain anchoring, direction exploration, prob-  
139 lem validation, method design, experiment execution, results integration, paper writing) with explicit return  
140 paths for backtracking, each with defined entry conditions, deliverables, and exit criteria (Fig. 2). Temporal  
141 ordering is enforced: experiments cannot begin before the evaluation protocol is frozen, and paper writing  
142 cannot begin before results are integrated.

143 This phase-gated structure addresses what we observe as the most consequential failure mode when AI  
144 agents conduct research without procedural constraints. In the protocol-free baseline of our controlled study  
145 (Project 6), the agent produced a complete manuscript while explicitly interleaving drafting with ongoing  
146 data extraction and analysis—beginning to write and compile the paper before the full analysis had finished,  
147 and then revising the draft as additional results arrived. In this baseline condition, intermediate deliverables  
148 (e.g., a locked protocol, sufficiency criteria, or an integration blueprint) were not externalized, reducing  
149 auditability and making goalpost-moving and stale-claim reuse difficult to detect. With the procedural work-  
150 flow imposed, the protocol requires explicit literature grounding, direction validation, evaluation-protocol  
151 locking before experimentation, and iteration/backtracking when evidence is insufficient—behaviours driven  
152 by the constraints rather than by additional model capability.

153 Among the phases of the research cycle, question formulation (direction exploration and problem vali-  
154 dation) proved both the most critical and the most difficult to automate well. In our projects, a single-agent  
155 literature review tended to produce conservative, incremental questions—extensions of existing work that  
156 minimize risk but also minimize novelty. We address this by combining structured multi-agent deliberation  
157 during Phase 1 with six explicit ideation strategies (contradiction mining, assumption challenging, cross-  
158 domain transfer, limitation-to-opportunity conversion, counterfactual reasoning, and trend extrapolation). In  
159 Project 3 (DESI), multi-agent brainstorming merged two initially separate directions into a more ambitious  
160 unified proposal; in Project 2 (ArchaicPainter), the divergent phase generated multiple candidate directions  
161 before convergence on the Li & Stephens extension that proved most productive.

## 162 Integrity discipline makes integrity checks explicit and auditable

163 The integrity discipline comprises seven persistent constraints that activate at defined trigger points and re-  
164 main active until project completion. Each targets a documented threat to research validity, but critically,  
165 each is enforced through explicit artifacts and halting/backtracking behaviour rather than implicit “best ef-  
166 fort” diligence (Table 1).

Table 1: Integrity discipline constraints and the failure modes they address. The final column summarizes examples of enforcement drawn from project artifacts and logs.

Constraint	Failure mode addressed	Observed enforcement
Metric immutability	Post hoc evaluation changes	Protocol locked at G2; changes require explicit authorization.
Complete reporting	Selective reporting	Negative results retained and disclosed (gpuADMIX; ArchaicPainter positive-only emission ablation).
Claim-evidence alignment	Unsupported assertions	Paper text reconciled with code/results; mismatches corrected.
Alternative-hypothesis exclusion	Confounded causal claims	Alternatives enumerated; discriminating tests required (DESI).
Reproducibility	Unreproducible computation	Seeds + environment/log capture + scriptable reruns required.
Verification	Unverified status claims	Fresh checks required; references verified; re-compute/re-compile before claims.
Figure standards	Substandard visualizations	Readable, consistent figures with explicit inclusion/caption checks before final compilation.

167 Across the protocol-constrained projects, these constraints were triggered repeatedly and resulted in  
168 concrete interventions: negative results were retained and disclosed (e.g., ArchaicPainter’s failed “positive-  
169 only” emission variant and gpuADMIX’s documented bugs/limitations), skeptical checks forced explicit  
170 alternative-hypothesis testing (DESI), reference and manuscript audits corrected unsupported citations and  
171 code–text mismatches, and verification steps enforced fresh checks (including reference verification and final  
172 compilation sanity checks) before claims were finalized.

## 173 Governance prevents indefinite continuation of failing approaches

174 The governance layer encodes four strategic functions: novelty assessment, scope control, failure manage-  
175 ment, and standards alignment. In our projects, governance repeatedly forced strategic self-assessment to  
176 be made explicit and auditable: novelty checks triggered redesigns when an idea was not yet publishable;  
177 scope control recorded and justified exclusions; and standards-alignment checks translated evidence suffi-  
178 ciency into concrete venue decisions and documented limitations when baselines or supplements could not

179 be completed.

180 The failure-management function proved particularly important: it enforces structured reassessment af-  
181 ter repeated setbacks and presents explicit pivot/downgrade/stop options, rather than allowing unbounded  
182 iteration to remain implicit. This mechanism mirrors the role of thesis-committee reviews and lab-group  
183 critiques, but externalizes the decision in a form that can be reviewed after the fact.

## 184 **Validation gates provide structured checkpoints**

185 The formalization includes four gates between major phases (G1–G4), each defining explicit criteria that  
186 must be satisfied before progression (Fig. 1). **G1** validates research direction, scientific significance, and  
187 resource feasibility. **G2** locks the evaluation protocol. **G3** confirms data quality and resource readiness. **G4**  
188 verifies that evidence is sufficient for writing.

189 These gates function as quality checkpoints within the methodology. In our validation projects, a human  
190 researcher reviewed gate criteria and provided approval or redirection—and this proved valuable, particu-  
191 larly when the researcher challenged an insufficiently novel plan at G2 (see Vignette 2). However, the gates  
192 are architecturally part of the methodology’s procedural logic, not a separate human-dependence mecha-  
193 nism: the criteria they enforce (e.g., “evaluation protocol must be locked before experimentation begins”)  
194 are themselves formalized constraints. The primary contribution is the methodology that structures what  
195 happens between gates—the phases and their explicit return paths, the integrity discipline, the governance  
196 interventions—which transforms AI behaviour regardless of who or what evaluates the gate criteria. Human  
197 participation at gates enhances quality by contributing scientific judgment that is difficult to formalize, but  
198 the methodology provides the research discipline that makes gate evaluation meaningful in the first place.

## 199 **Methodology formalization adapts across research types**

200 A single set of methodological principles must accommodate the genuine diversity of research practice to  
201 be considered a valid formalization. Performance-driven method development requires different procedures  
202 from story-driven scientific discovery or utility-driven tool evaluation. We accommodated this diversity  
203 through differential activation: the integrity discipline and governance layers operate identically across all  
204 research types (metric immutability and complete reporting are universal), while the procedural workflow  
205 adapts. **Method** projects emphasize evaluation protocol locking, baseline reproduction, and mandatory it-  
206 eration (minimum three diagnose–hypothesize–fix–measure cycles). **Discovery** projects emphasize analy-

207 sis storyboard design, alternative hypothesis exclusion, and narrative coherence. **Tool** projects emphasize  
208 benchmarking, usability, and documentation. **Hybrid** projects activate both Method and Discovery tracks.

209 This structure suggests an organizing principle: research methodology has a *universal core* of integrity  
210 and governance norms, with *type-specific procedural instantiations* that adapt to the nature of the knowledge  
211 claim being made.

## 212 **Structured multi-agent deliberation provides repeated, role-differentiated review checkpoints**

213 In human research, critical decisions are repeatedly stress-tested through discussions at multiple points in  
214 the workflow (lab meetings, advisor feedback, collaborator critique), not only at the journal-review stage.  
215 We operationalize an analogous mechanism by embedding structured multi-agent deliberation at five criti-  
216 cal junctures. At each juncture, the system dynamically instantiates role-specialized *sub-agents* conditioned  
217 on the research question and current project state (e.g., a domain expert, a skeptical critic, and an editor).  
218 Crucially, each sub-agent is given an independent context and evaluates the same artifact (protocol, plan,  
219 results integration blueprint, or manuscript section) against a shared rubric without sharing its intermediate  
220 reasoning with the others. Deliberation proceeds in structured rounds (maximum five): convergence re-  
221 quires unanimous PASS; otherwise, the artifact is modified and *all* sub-agents re-assess the full artifact. The  
222 resulting recommendations can trigger explicit returns to earlier phases (e.g., rerun experiments, redesign  
223 analyses, or retract unsupported claims). Unresolved disagreements are escalated to the human researcher.

224 This draws on evidence that diverse multi-agent panels improve reasoning quality [31–33] and extends  
225 the principle from single-query tasks to sustained, multi-phase research where accumulated context matters.

## 226 **Cross-domain validation**

227 We applied the complete methodology formalization to research projects spanning diverse scientific do-  
228 mains and research types, without any domain-specific modification (Supplementary Papers; Fig. 5). In  
229 each protocol-constrained project, a general-purpose LLM operating under the full constraints conducted the  
230 complete research lifecycle. We report here on five protocol-constrained projects spanning population ge-  
231 nomics, computational bioinformatics, human evolutionary genetics, computational population genetics, and  
232 condensed-matter physics. We additionally report a sixth project (Project 6): a controlled study on the same  
233 dataset and task, in which the same AI agent produced two complete manuscripts under matched conditions  
234 with and without the protocol (Supplementary Papers 6A/6B).

235 Across these projects, we focus on a small set of evidence-backed research outputs that are directly au-  
 236 ditable in the generated papers: (i) closed-form analytical derivations that can be independently checked,  
 237 (ii) quantitative ablations that resolve specific modeling design choices, (iii) algorithmic refactors that pre-  
 238 serve the objective but change the computational primitive, and (iv) domain analyses that recapitulate known  
 239 signals while surfacing limitations and inconsistencies transparently. Table 2 summarizes these outputs and  
 240 provides an evidence anchor (equation/figure/table) for each.

Table 2: **Evidence-backed cross-domain research outputs produced under the protocol.** Each entry is stated narrowly and paired with an explicit evidence anchor in the corresponding Supplementary Paper.

Project	Output type	Narrow, checkable claim	Evidence anchor
gpuADMIX	Algorithmic refactor	The binomial-likelihood EM updates can be rearranged into a small fixed set of GEMM (matrix multiplications) plus elementwise operations, making the original objective GPU-native without changing the model.	Supp. Paper 1
ArchaicPainter	Ablation conclusion	In introgression HMM emission, treating mismatches as informative (bidirectional emission) is necessary for segment discriminability when the reference matches the donor; down-weighting mismatches reduces simulation F1 from 0.480 to 0.154.	Supp. Paper 2
HapGraph	Analytical correction	For cross-population IBD between two modern lineages, the expected mean segment length scales as $50/T$ cM (not $100/T$ ); with detection threshold $\ell_{\min}$ , $E[L   L > \ell_{\min}] = 50/T + \ell_{\min}$ .	Supp. Paper 4
z2-quantum-mpemba	Closed-form derivation	For product initial states, the initial $Z_2$ entanglement asymmetry admits a closed-form expression bounded by $\ln 2$ , providing an exact calibration scale for numerical experiments.	Supp. Paper 5
DESI	Robust observation (+ caveat)	Window-level mean pairwise TMRCA is substantially deeper within AFR (1,229 ka) than EAS (859 ka) and EUR (920 ka) across all 22 autosomes; a FitCoal-parameter stress test suggests mismatch under a limited model class and summary statistics, motivating follow-up with better-matched simulation and likelihood-based fitting.	Supp. Paper 3

241 **Project 1 (gpuADMIX, Type H—Method + Tool):** GPU-accelerated ancestry estimation preserving  
 242 the exact ADMIXTURE binomial likelihood. A key method insight is that the EM updates can be alge-  
 243 braically reorganized into a small fixed set of GEMM operations plus elementwise steps, so the original  
 244 binomial objective becomes GPU-native without changing the model. The constrained agent developed a  
 245 complete software tool (Python/PyTorch), conducted ablation experiments across five random seeds at each  
 246 of nine  $K$  values, and produced a publication-ready manuscript. The procedural workflow enforced evalua-  
 247 tion protocol locking before experiments, and the integrity discipline required reporting all seeds including  
 248 failures. A critical integrity event occurred when stored results for  $K \geq 8$  were discovered to originate  
 249 from an earlier code version; the verification constraint forced re-measurement, which reversed the paper’s

250 narrative from “limited at high  $K$ ” to “competitive at all  $K$  via multi-seed strategy.” Multi-agent delibera-  
251 tion panels caught two fatal citation errors (one entirely fabricated reference), a formula mismatch between  
252 code and text (FISTA momentum vs. simple Nesterov), and a misattributed benchmark figure—all before the  
253 manuscript was finalized. The resulting tool achieved a  $213\times$  speedup over ADMIXTURE and  $41\times$  over  
254 fastmixture while maintaining  $Q$ -matrix correlation with ADMIXTURE outputs  $r^2 > 0.9999$ .

255 **Project 2 (ArchaicPainter, Type H—Method + Discovery):** A three-state Hidden Markov Model ex-  
256 tending the Li & Stephens haplotype-matching framework to archaic reference genomes for per-haplotype  
257 introgression detection. The constrained agent developed the method, validated it on 50 independent simula-  
258 tion replicates ( $F_1 = 0.480 \pm 0.095$ ;  $56\times$  improvement over a density baseline, Wilcoxon  $p = 8.9 \times 10^{-16}$ ),  
259 and applied it to chromosome 21 of the 1000 Genomes Project. A central design conclusion is that mis-  
260 matches to the archaic reference are not mere noise but carry discriminative signal when the reference is a  
261 faithful proxy for the donor: down-weighting mismatches via a positive-only emission reduces simulation  
262  $F_1$  to  $0.154 \pm 0.079$ . Multi-agent deliberation identified three fatal and four major evidence gaps, forcing  
263 the agent back from Phase 5 to Phase 4 for supplemental experiments. The human researcher overrode the  
264 requirement for direct tool comparison (HMMix/DAIseg installation failed), and the governance layer en-  
265 sured this was documented as an acknowledged limitation with a corresponding venue downgrade. Notably,  
266 the method independently recovered  $\sim 2\%$  Neanderthal ancestry in European and East Asian populations—  
267 consistent with published estimates from independent methods [34]—and detected the complete interferon  
268 receptor gene cluster (*IFNAR1*, *IFNAR2*, *IFNGR2*, *IL10RB*) as a Neanderthal-introgressed region in Euro-  
269 peans, corroborating published adaptive introgression findings [35].

270 **Project 3 (DESI, Type D—Discovery):** A window-based coalescent-depth analysis quantifying genome-  
271 wide mean pairwise TMRCA differences across super-populations. This project exhibited the most conse-  
272 quential governance interventions. At G2, the human researcher challenged the proposed plan (“is there  
273 actually a new method here, or are you just running existing tools?”), triggering a complete project redesign.  
274 During execution, the AI agent retracted its own interpretation of initial results, explicitly stating “my ear-  
275 lier interpretation was wrong” and “is this enough for *Nature Genetics*? No—not from this angle.” The  
276 resulting paper grounds its claims through simulation calibration and chromosome-wide consistency checks:  
277 within-AFR mean TMRCA is  $1,229.5 \pm 6.4$  ka, deeper than within-EAS ( $858.8 \pm 5.4$  ka) and within-EUR  
278 ( $920.2 \pm 5.5$  ka), with the AFR–EAS difference ( $370.7 \pm 8.4$  ka) positive on all 22 autosomes. The paper  
279 further proposes—based on a parameter scan of a *panmictic*  $\sim 930$  ka bottleneck model and window-level

280 summary statistics (e.g.,  $P(\bar{T}_{\text{AFR}} > 930 \text{ ka})$  and  $\bar{T}_{\text{AFR}}$ )—that the bottleneck may be less severe than Fit-  
281 Coal’s published parameterization. However, after manual review we note that this inference depends on the  
282 specific summary-statistic comparison and simulation setup, and that the simulated data use a much smaller  
283 number of windows than the real-data analysis, which may bias the tail-fraction estimate; we therefore treat  
284 this result as a suggestive direction for follow-up with better-matched simulation and likelihood-based fitting  
285 rather than a decisive conclusion.

286 **Project 4 (HapGraph, Type C—Tool):** A Bayesian admixture graph inference tool unifying F-statistics  
287 and inter-population IBD sharing for joint estimation of graph topology and admixture proportions ( $\alpha$ ), with  
288 an IBD-based timing model for  $T$  on fixed topologies. This project demonstrated the integrity discipline’s  
289 role in iterative bug discovery and correction. During Phase 4a prototype testing, the verification constraint  
290 identified that the  $F_2$  path-length likelihood approximation was insensitive to admixture—a critical imple-  
291 mentation flaw that rendered the tool unable to estimate  $\alpha$ . Rather than proceeding with a partially working  
292 system, the constraint forced a complete rewrite to an ancestry-vector formulation with NNLS branch-length  
293 fitting, followed by systematic correction of two additional issues ( $F_3$  bias correction per Patterson et al.  
294 2012, and IBD source-clade filtering) through three diagnose–fix–measure iterations in Phase 4b. A techni-  
295 cally important correction in this project is the cross-population IBD length scale: for two modern lineages  
296 separated by  $T$  generations, the expected mean length is  $50/T$  cM, and under a detection threshold  $\ell_{\min}$   
297 it shifts to  $50/T + \ell_{\min}$ . Multi-agent deliberation on the manuscript draft then caught critical discrepan-  
298 cies between the paper text and the actual code and forced reconciliation before the manuscript could pass  
299 the quality gate. On simulation benchmarks (S1–S3, 20 seeds each), HapGraph achieved 100% topology  
300 accuracy at  $K \leq 2$  (where  $K$  denotes the number of admixture edges),  $\alpha$  MAE of 0.054, and  $T$  MAE  
301 of 10.9 generations with  $T$  95% CI coverage of 93%. Applied to 26 populations from the 1000 Genomes  
302 Project, the tool identified  $K = 8$  admixture events without manual curation; seven of the eight targets  
303 correspond to well-studied admixed populations (ASW, ACB, PUR, CLM, MXL, BEB, PJL), and timing  
304 estimation for real data was deferred pending whole-genome IBD pre-computation.

305 **Project 5 (z2-quantum-mpemba, Type D—Discovery):** A condensed-matter physics study of the  $Z_2$   
306 quantum Mpemba effect in integrable one-dimensional free-fermion chains (transverse-field Ising model and  
307 anisotropic XY chain). The constrained agent derived an exact closed-form expression for the initial entan-  
308 glement asymmetry and verified it to machine precision, then used exact diagonalization to establish that  
309  $Z_2$  QME is present across all 23 post-quench field values tested—including the DQPT-free ferromagnetic

310 phase. Statistical analysis of 671 initial-state pairs revealed a strong correlation between the Mpemba cross-  
311 ing time  $t_M$  and the initial entanglement-asymmetry imbalance (Spearman  $\rho = +0.90$ ,  $p \approx 3 \times 10^{-240}$ ),  
312 while showing that DQPT singularities modulate but do not govern crossings (KS  $p \approx 6 \times 10^{-25}$ ; proximity  
313 test  $p = 0.27$ ). Finite-size scaling over  $N = 8$ –18 supported persistence in the thermodynamic limit. No-  
314 tably, despite modest computational runtime (order tens of minutes on a 128-core CPU server), this project  
315 required sustained analytical reasoning and exact-solution verification, consistent with the effort distribution  
316 in human theoretical physics.

317 Across the five protocol-constrained projects (Supplementary Table S1), the procedural workflow com-  
318 pleted all seven phases with all four gates enforced in each. The integrity discipline repeatedly halted pro-  
319 gression until verification and documentation requirements were satisfied (Table 1), and multi-agent de-  
320 liberation served as an embedded adversarial review mechanism that surfaced actionable gaps (evidence  
321 insufficiency, code–text mismatches, and citation integrity issues) before manuscripts were finalized. Gov-  
322 ernance interventions—including methodology pivots, scope reductions, and venue-alignment downgrades  
323 when evidence was insufficient—were required in multiple projects, illustrating that sustained strategic self-  
324 assessment can be operationalized as part of the protocol rather than left to ad hoc prompting.

### 325 **Controlled study (Project 6): protocol increases auditability and methodological structure**

326 To isolate the effect of the protocol from underlying model capability, we report a matched controlled study  
327 (Project 6) in which the sole manipulated factor is activation of the formalized constraints. The same base  
328 model (Claude Opus 4.6) operating in the same IDE (Cursor) was tasked with conducting autonomous  
329 scientific research on the 1000 Genomes Project 2022 high-coverage dataset and producing a complete  
330 manuscript. In the protocol-free baseline condition, the agent produced a technically solid, indel-focused  
331 analysis (target standard: *Genomics*) with a single-file manuscript and a set of analysis scripts. In the Cur-  
332 sor+Amplify condition, the agent produced a modular manuscript and, crucially, externalized the method-  
333 ological process as auditable protocol artifacts (intake anchor, literature review and gap analysis, explicit  
334 planning documents, and an integration blueprint), enabling structured iteration and clearer evidence-to-  
335 claim traceability (Table 3). Both complete manuscripts (with protocol vs. protocol-free) are included as  
336 Supplementary Papers 6A and 6B.

337 To make the difference in artifact externalization concrete, Fig. 3 summarizes the produced directory  
338 structure in each condition.

Table 3: **Controlled study: artifact-level comparison of protocol-free baseline versus Cursor+Amplify.** Both conditions used the same base model and tool access on the same dataset; differences are assessed via the presence of auditable artifacts and executable outputs.

Audit target	Cursor+Amplify protocol	Protocol-free baseline
Protocol deliverables externalized	Intake + literature + plan + integration artifacts present (research anchor; literature review; gap analysis; sufficiency and alternative explanations; argument blueprint)	No protocol deliverables directory; analysis proceeds ad hoc
Manuscript organization	Modular LaTeX (preamble + 5 section files)	Single-file LaTeX manuscript
Executable analysis pipeline	9 scripts under scripts/; figures and paper maintained as separate top-level outputs	13 scripts under analysis/scripts; paper maintained under analysis/manuscript
Figures integrated into manuscript	4 includegraphics blocks in main text	5 includegraphics blocks in main text
Analysis scope (as implemented)	Unified cross-type population-structure comparison using SNVs, INDELS, and SVs	Indel mutation-spectrum analysis with NMF decomposition

### 339 Mechanism: how constraints alter behaviour in practice

340 To illustrate *how* constraints change agent behaviour—not merely that outcomes differ—we trace three rep-  
 341 resentative episodes from the validation projects.

342 **Vignette 1: Verification constraint catches stale evidence (gpuADMIX).** During Phase 4, the agent  
 343 benchmarked gpuADMIX against fastmixture at  $K = 2-10$  and initially concluded that gpuADMIX was  
 344 inferior at  $K \geq 8$ . The agent flagged this honestly (anti-cherry-pick constraint: “this needs honest dis-  
 345 closure”). When the human researcher later questioned the claim, the verification constraint required fresh  
 346 re-measurement rather than reliance on stored results. Re-running revealed that the  $K \geq 8$  numbers origi-  
 347 nated from an earlier, buggy code version: the true best-of-five gpuADMIX log-likelihood matched or ex-  
 348 ceeded fastmixture at every  $K$ . Without the verification constraint, the stale results would have propagated  
 349 into the manuscript unchallenged. Additionally, multi-agent review of the paper draft caught a completely  
 350 fabricated reference (a non-existent paper attributed to a real author) and a formula discrepancy between  
 351 the code (FISTA momentum) and the text (simple Nesterov)—errors that would have survived conventional  
 352 self-review.

(a) Cursor+Amplify protocol	(b) Protocol-free baseline
<p><b>docs/</b> (protocol deliverables) 01_intake/research-anchor.yaml 02_literature/ (review + gaps) 03_plan/ (locked plan) 05_integration/ (claim–evidence)</p> <p><b>paper/</b> (modular LaTeX) main.tex + preamble.tex sections/ (5 section files) supplementary/</p> <p><b>scripts/</b> (9 analysis scripts) <b>figures/</b> (plots + exports)</p>	<p><b>analysis/</b> <b>manuscript/</b> main.tex (single file) <b>scripts/</b> (13 analysis scripts) <b>figures/</b> (plots)</p> <p><i>No explicit protocol deliverables directory (intake/literature/plan/integration artifacts not externalized)</i></p>

Figure 3: **Controlled study: artifact externalization differs under a phase-gated protocol.** In the Cursor+Amplify condition, protocol deliverables are externalized as auditable on-disk artifacts (docs/) alongside a modular manuscript (paper/) and an executable analysis pipeline (scripts/). In the protocol-free baseline, analysis proceeds without an explicit protocol artifact directory and outputs are organized around scripts and a single-file manuscript.

353 **Vignette 2: Governance forces project redesign; integrity forces claim retraction (DESI).** At the G2  
354 gate, the human researcher challenged the proposed analysis plan: “Is there actually a new method here,  
355 or are you just running existing tools?” The agent honestly acknowledged that the plan was “an analysis  
356 pipeline plus new application, not a new method,” triggering a complete project redesign. Later, during  
357 Phase 4, the agent’s initial finding—a two-component TMRCA distribution suggesting ancient structure—  
358 was questioned by the human. The agent re-analysed the evidence, explicitly stated “my earlier interpre-  
359 tation was wrong,” and downgraded the claim: “Is this enough for *Nature Genetics*? No—not from this  
360 angle.” Rather than defending the original interpretation, the integrity discipline required the agent to pur-  
361 sue additional simulation experiments and calibration checks. In the final write-up, a technical bias check  
362 with random labels (Model A) yields an AFR–EAS difference of  $0.5 \pm 0.3$  ka, while a realistic Out-of-  
363 Africa model with panmictic deep ancestry (Model OoA) predicts an AFR–EAS difference of 463.9 ka.  
364 The empirical window-level AFR–EAS difference (370.7 ka) is then assessed in the context of model-  
365 based stress tests, including a FitCoal-parameter scan that compares window-level summary statistics (e.g.,  
366  $P(\bar{T}_{AFR} > 930 \text{ ka})$  and  $\bar{T}_{AFR}$ ) from 600 simulated windows against 27,507 real-data windows, and there-  
367 fore explicitly scopes conclusions to the tested model class and summary-statistic estimator. The governance  
368 layer’s failure-management function—forcing structured reassessment after repeated setbacks—was essen-  
369 tial to reaching this outcome rather than abandoning the project or advancing unsupported claims.

370 **Vignette 3: Multi-agent deliberation catches code–text mismatches (HapGraph).** During Phase 6 (pa-  
371 per writing), three AI review agents independently audited the HapGraph Methods section against the actual  
372 codebase. The adversarial reviewer flagged multiple code–text inconsistencies that would survive conven-  
373 tional self-review (the formulas were plausible and internally consistent, and compilation/unit tests would  
374 not catch them). One representative example was the inter-population IBD timing model: the corrected im-  
375 plementation and write-up require  $\mathbb{E}[\bar{L} \mid \bar{L} > \ell_{\min}] = 50/T + \ell_{\min}$ , reflecting a breakpoint rate of  $1/(2T)$   
376 per Morgan for two modern lineages. The claim-evidence alignment constraint required the writing agent  
377 to reconcile every quantitative assertion with the implementation before the section could pass the qual-  
378 ity gate, and the multi-round deliberation protocol ensured discrepancies were resolved before the text was  
379 finalized. Additionally, the same review cycle identified fabricated bibliography metadata (correct author  
380 names paired with wrong titles, volumes, and DOIs) that would have been undetectable without systematic  
381 cross-referencing.

### 382 **Relationship to existing approaches**

383 The methodology formalization addresses a dimension distinct from the capabilities that existing AI research  
384 systems optimize. To clarify this distinction, we categorize recent systems by what they encode and what  
385 they omit (Table 4).

386 **Domain-specific task automation.** Systems such as AutoBA [20], BioMaster [21], and Biomni [14] en-  
387 code knowledge of *which tools to run and in what order* for specific scientific domains. AutoBA and  
388 BioMaster automate bioinformatics pipelines (RNA-seq, ChIP-seq, spatial transcriptomics) by generating  
389 and executing analysis plans from minimal user input; Biomni extends this to 25 biomedical subfields by  
390 orchestrating 105 software tools, 150 biological protocols, and 59 databases. Similarly, ChemCrow [36]  
391 and SciAgents [37] achieve deep integration with chemistry and materials science knowledge, respectively.  
392 These systems excel at executing established workflows—a genuine contribution—but do not encode higher-  
393 level methodological reasoning: they do not question whether a planned analysis adequately addresses the  
394 research question, require evaluation protocols to be locked before execution, or enforce complete reporting  
395 of unfavourable outcomes.

396 **End-to-end autonomous AI scientists.** A rapidly growing category of systems aims to automate the full  
397 research cycle. The AI Scientist [15] and its successor AI Scientist v2 [16] generate ideas, write code, run  
398 experiments, and produce complete ML papers—the latter achieving the first AI-generated peer-reviewed  
399 publication at an ICLR 2025 workshop. Kosmos [17] (FutureHouse) performs extended 12-hour discovery  
400 sessions with a structured world model that maintains coherence across ~200 agent rollouts, producing re-  
401 ports with 79.4% statement accuracy across metabolomics, materials science, neuroscience, and statistical  
402 genetics. Robin [18] (FutureHouse) achieved the first fully automated biological discovery—identifying a  
403 drug candidate for age-related macular degeneration through iterative hypothesis generation and experimen-  
404 tal design. Google’s AI-powered empirical software system [38] uses tree search to optimize scientific code  
405 across six benchmarks, achieving expert-level performance in genomics, epidemiology, and neuroscience.  
406 Agentomics [39] autonomously develops state-of-the-art ML models for biomedical datasets, outperforming  
407 human expert solutions on 11 of 20 benchmarks. These systems demonstrate that LLMs can execute substan-  
408 tive research tasks autonomously—a capability our approach leverages but does not aim to replicate. What  
409 they do not provide is persistent methodological discipline: none enforces evaluation immutability, requires  
410 complete reporting of negative results, mandates alternative-hypothesis exclusion before causal claims, or  
411 provides governance mechanisms to halt failing approaches.

412 **Structured hypothesis generation.** Google’s AI co-scientist [19] is the closest existing system to our  
413 approach in terms of structured process. Built on Gemini 2.0, it employs specialized agents (Generation,  
414 Reflection, Ranking, Evolution, Proximity, Meta-review) in a “generate, debate, and evolve” cycle that it-  
415 eratively improves research hypotheses through tournament evolution, with demonstrated validation in drug  
416 repurposing, novel target discovery, and bacterial evolution. The system shares our philosophy that multi-  
417 agent deliberation and human-AI collaboration improve research quality. However, the two approaches  
418 differ fundamentally in *scope* and *what is formalized*. AI co-scientist formalizes the hypothesis genera-  
419 tion and refinement stage—producing ranked candidate hypotheses for scientists to validate experimentally.  
420 Our formalization covers the *complete research lifecycle*: from question formulation through experiment  
421 execution, results integration, and manuscript preparation, with persistent integrity constraints (evaluation  
422 immutability, complete reporting, claim-evidence alignment) active throughout. AI co-scientist does not en-  
423 code constraints on how experiments should be conducted once a hypothesis is selected, does not require  
424 complete reporting of negative outcomes, and does not provide governance mechanisms to force pivoting

425 when approaches fail. The relationship is again complementary: AI co-scientist’s sophisticated hypothesis  
 426 generation could feed into our methodology’s Phase 1 (direction exploration), after which our integrity and  
 427 governance layers would ensure the selected hypothesis is pursued with sustained methodological discipline.

428 **Knowledge synthesis and deep research.** Products such as OpenAI’s Deep Research [22] and Google’s  
 429 Gemini Deep Research automate literature search and report generation at scale. These are valuable for the  
 430 direction exploration phase of the research cycle but cover only one component of the full methodology—  
 431 they do not design experiments, lock evaluation protocols, or verify that claims rest on specific evidence.

432 **Complementarity.** Our formalization is complementary to all these approaches rather than competitive  
 433 with them. Domain-specific systems encode *what tools to use*; autonomous scientists encode *the capabil-*  
 434 *ity to execute*; deep research products encode *how to synthesize existing knowledge*; our work encodes *the*  
 435 *methodological discipline that makes research outputs reliable*. In principle, the three constraint layers (pro-  
 436 cedural workflow, integrity discipline, governance) could be applied atop any of these systems—providing  
 437 the methodological scaffolding that currently separates impressive capability demonstrations from publish-  
 438 able, reproducible science.

Table 4: Methodology formalization addresses a dimension distinct from domain capability. Each column represents a category of existing systems; rows indicate whether the category provides each capability. Entries reflect capabilities as reported in the cited system descriptions. Legend: ✓=core feature; Partial=limited or not systematic; —=not reported/primary goal.

	<b>This work</b>	<b>Autonomous scientists</b>	<b>AI co-scientist</b>	<b>Domain agents</b>	<b>Deep re-search</b>
Persistent integrity	✓	—	—	—	—
Eval. immutability	✓	—	—	—	—
Complete reporting	✓	—	—	—	—
Human approval gates	✓	Partial	Partial	Partial	Partial
Direction ideation	✓	✓	✓	Partial	Partial
Structured deliberation	✓	Partial	✓	—	—
Autonomous execution	✓*	✓	✓	✓	✓
Tool execution	✓ <sup>†</sup>	Partial	Partial	✓	Partial
Domain agnostic	✓	Partial	✓	—	✓

\* Amplify is designed so that phase transitions are subject to explicit human review/approval at validation gates. When desired, a simple prompt setting can enable a “fast mode” in which the agent proceeds by default using the protocol’s recommended decisions, minimizing human intervention while retaining the same phase structure and constraints.

<sup>†</sup> Domain tool access is provided by the execution environment (e.g., Cursor-style IDEs that support web search/browsing and a bash shell), allowing agents to install and invoke many domain tools on demand.

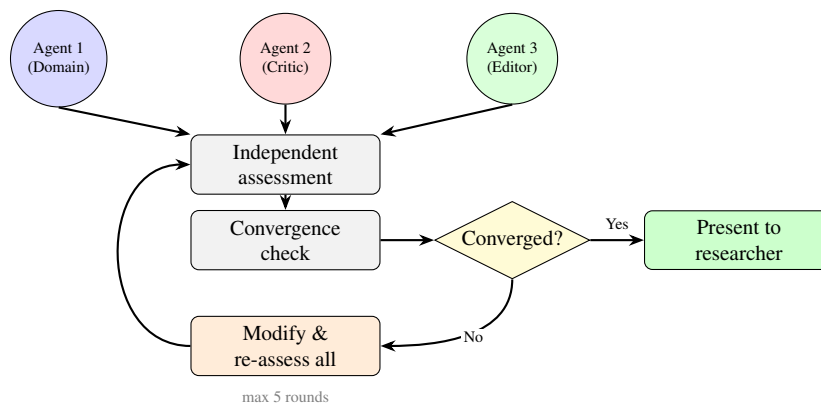


Figure 4: **Multi-agent deliberation protocol.** Three role-specialized sub-agents are dynamically instantiated based on the research question and current project state. Each sub-agent assesses the same artifact with an independent context before any convergence check, reducing anchoring and enabling adversarial critique. Deliberation proceeds through structured rounds: convergence requires unanimous PASS; otherwise the artifact is modified and all sub-agents re-assess the full artifact. Unresolved disagreements are escalated to the human researcher.

## 439 Discussion

440 The most surprising aspect of these results is not that AI agents can be constrained to follow rules—  
441 instruction-following is a known capability of modern LLMs—but that a relatively compact formalization of  
442 research methodology suffices to transform the behaviour of general-purpose AI from ad hoc text generation  
443 into structured research practice across diverse domains. Two findings stand out. First, the complete research  
444 process—how to formulate questions, validate directions, design methods, iterate through failures, recognize  
445 when to pivot, integrate results into evidence-grounded narratives—*can* be articulated as a coherent system  
446 of interlocking constraints. Second, AI agents, when given this system, can be steered to follow it: in our  
447 validation projects, the same models that without constraints can produce complete manuscripts while leav-  
448 ing planning and verification implicit instead conduct multi-phase, methodologically disciplined research  
449 with auditable intermediate artifacts. This suggests that the gap between “AI that can write about science”  
450 and “AI that can do science responsibly” is not primarily a gap in model capability but a gap in *methodology*  
451 *transfer*: the procedural knowledge that experienced researchers possess was simply never given to the AI.

452 This finding resonates with a broader insight from constitutional AI [40], which demonstrated that ex-  
453 plicit principles steer model behaviour more reliably than implicit training signals. We extend this from AI  
454 safety to scientific integrity, and from individual interactions to multi-phase projects with persistent state  
455 and accumulated constraints. The three-layer decomposition we propose (workflow, discipline, governance)

**End-to-end validation set used in this paper (six projects).**

**Projects 1–5 (protocol-constrained; Cursor+Amplify):**

**gpuADMIX** (population genetics tool/method; target: Bioinformatics)

**ArchaicPainter** (paleogenomics method+discovery; target: Genome Biology)

**DESI** (deep-time demography discovery; target: Nature Genetics)

**HapGraph** (admixture-graph tool; target: Bioinformatics)

**z2-quantum-mpemba** (condensed-matter discovery; target: PRB/SciPost)

**Project 6 (controlled study; same task, two complete papers):**

**Paper 6A** (with protocol) vs **Paper 6B** (protocol-free baseline)

Human genomics on 1000GP (INDEL-focused analysis; target: Genomics)

Figure 5: **Cross-domain validation set and controlled study.** Five projects were executed under the full protocol constraints (Cursor+Amplify) across population genetics, paleogenomics, human evolutionary genetics, computational population genetics, and condensed-matter physics. A sixth project is a controlled study on the same 1000 Genomes dataset and task, producing two complete papers under matched conditions with (Paper 6A) and without (Paper 6B) the protocol.

456 mirrors how methodology functions in human research groups—the PI’s plan, the lab’s integrity culture, the  
457 oversight of committees [7]—and the fact that each layer catches distinct failure modes (see Results Section)  
458 suggests that these are not arbitrary categories but reflect genuine functional divisions in how reliable knowl-  
459 edge is produced. Methodological constraints do not replace domain expertise, stronger models, or improved  
460 tool integrations. Rather, they make a research agent’s methodological *process* explicit and auditable—  
461 externalizing phase progression, turning integrity checks into enforceable obligations with concrete artifacts,  
462 and forcing stop/pivot decisions to be written down and justified. A practical implication is that the formal-  
463 ization can be applied directly atop general-purpose LLMs. Because the protocol is model-agnostic, it can  
464 benefit automatically from improvements in underlying foundation models (knowledge, reasoning, tool use)  
465 while keeping the methodology layer fixed and auditable. In cost terms, this shifts effort from large-scale  
466 model training to an executable, reusable protocol that can be shared and iterated as open-source software.

467 The validation projects also reveal a domain-dependent shift in where effort is required. In data-analysis  
468 projects, the dominant cost is often computational execution and debugging of pipelines. In contrast, in the  
469 z2-quantum-mpemba condensed-matter project, computational runtime was modest but the agent devoted  
470 substantially more effort to analytical reasoning and proof-like verification against exact constraints (closed-  
471 form initial conditions, finite-size scaling, and benchmarkable dynamical signatures)—a pattern that aligns  
472 with how human theoretical physics often emphasizes reasoning over data processing.

473 Looking forward, one can view protocol-layer enforcement as a stepping stone rather than an endpoint.  
474 In principle, the same constraints could be used as training signals (e.g., via reinforcement learning or pref-

475 erence optimization) so that models internalize parts of the methodology rather than relying on external  
476 enforcement; we did not pursue such training due to resource constraints and because our goal here is an  
477 auditable, model-agnostic protocol. We also observed meaningful performance differences across founda-  
478 tion models on different phases and tasks. A protocol controller makes it straightforward to exploit this  
479 heterogeneity by routing different phases or deliberation sub-agents to models that are better suited for the  
480 corresponding cognitive demands (e.g., long-horizon planning, code execution, or mathematical reasoning).

481 Several risks require candid acknowledgment. First, formalized methodology enforces *process* but can-  
482 not assess *substance*: a project can satisfy every constraint and still be scientifically trivial. Moreover, perfor-  
483 mance remains bounded by the underlying model’s knowledge and reasoning ability; the protocol can force  
484 checks and backtracks, but it cannot invent missing domain insight. Human judgment at gates addresses  
485 this partially, but users should understand that methodology-constrained AI produces rigorously *structured*  
486 research, not independently *validated* science. Second, the current formalization encodes the hypothetico-  
487 deductive methodology dominant in quantitative sciences; qualitative research, abductive reasoning, and  
488 indigenous knowledge systems follow different methodological traditions not captured here. Extending the  
489 formalization to accommodate methodological pluralism is essential future work. Third, constraints are en-  
490 forced through structured instructions rather than formal verification; LLMs are probabilistic systems, and  
491 subtle violations (such as biased interpretation of ambiguous results) may evade detection.

492 Overall, the results support a practical path to scaling methodological rigour: treat methodology as  
493 an auditable, reusable protocol layer that can be applied to general-purpose models today, improved as  
494 foundation models advance, and eventually partially internalized through training or deployed via phase-  
495 and role-specific model routing.

## 496 **Methods**

### 497 **Approach to methodology formalization**

498 To formalize scientific methodology, we analysed how experienced researchers maintain rigour across the  
499 research lifecycle and distilled recurring patterns into explicit, executable principles. Sources included pub-  
500 lished methodology guides, research integrity literature [29, 30, 41], reporting standards (CONSORT [3],  
501 PRISMA [4]), pre-registration frameworks [1, 2], and the accumulated conventions of peer review prac-  
502 tice. We sought principles that were (a) domain-agnostic—applicable to computational research regard-

503 less of field, (b) enforceable—expressible as verifiable constraints rather than aspirational guidelines, and  
504 (c) decomposable—assignable to a specific functional layer without entangling with other concerns.

505 This process yielded 24 distinct principles organized into the three-layer structure described in Results.  
506 Six overarching constitutional rules govern the entire formalization: (1) research type determines all proce-  
507 dural paths; (2) the target publication standards are established at inception and influence all downstream  
508 decisions; (3) the AI agent adopts a domain-specific expert identity rather than operating as a generic assis-  
509 tant; (4) methodological justification precedes implementation; (5) evaluation criteria, once locked, cannot be  
510 modified without explicit human authorization; and (6) no claim may be made without fresh computational  
511 verification.

## 512 **Implementation as a phase-gated research protocol**

513 To test whether formalized methodology can be transferred to AI agents as an operational protocol (rather  
514 than a single prompt), we implemented each principle as a structured protocol card specifying trigger condi-  
515 tions, required inputs, procedural steps, deliverables, exit criteria, and integration points. Collectively, these  
516 cards define a stateful, phase-gated process: phases advance only when entry/exit criteria are met; invariants  
517 such as metric immutability persist across phases; and violations trigger halts and explicit return paths for  
518 correction. Constraints are categorized as *rigid* (followed without adaptation; e.g., metric immutability, com-  
519 plete reporting) or *flexible* (principles adapted to context; e.g., direction exploration, method design). Rigid  
520 constraints include explicit verification steps—for example, the metric immutability constraint requires com-  
521 parison of any proposed evaluation change against the locked protocol, with mandatory halt on discrepancy.

522 The protocol operates atop general-purpose LLMs (Claude, GPT-4, Gemini, and others) without fine-  
523 tuning or model modification. At runtime, a lightweight protocol controller ensures that the agent is operating  
524 under the correct phase and discipline rules by loading the relevant protocol cards into the agent’s working  
525 context as structured system-level instructions. Crucially, persistent state is maintained outside the model  
526 via logged artifacts and explicit interfaces between phases, so that enforcement is not merely “remembering  
527 a prompt” but checking protocol invariants against concrete outputs (e.g., the locked evaluation protocol,  
528 recorded seeds, and claim–evidence tables).

529 **Controller implementation in the open-source release.** Amplify is implemented as a plugin-style skills  
530 library plus a small set of runtime hooks and always-on rules that make protocol loading explicit and repro-

531 ducible. The skills library is declared in a plugin manifest (`amplify/.cursor-plugin/plugin.json`) and stored  
532 as protocol cards under `amplify/skills/*/SKILL.md`. A `SessionStart` hook (`amplify/hooks/hooks.json` calling  
533 `amplify/hooks/session-start.sh`) injects the full content of the `using-amplify` skill into the system context at  
534 the beginning of a chat session, making global workflow constraints (e.g., one-phase-per-turn, gate enforce-  
535 ment, and mandatory skill invocation when triggers apply) available before any phase begins. For Cursor-  
536 style rule systems, an equivalent always-on bootstrap rule file (`amplify/install/amplify-bootstrap.mdc`) pro-  
537 vides the same controller constraints when hook-based plugins are not used.

538 Phase transitions are enforced procedurally: the protocol requires the agent to stop after completing a  
539 phase or reaching a gate, present deliverables and a gate checklist, and wait for explicit user approval before  
540 proceeding. Discipline and governance cards (e.g., `metric-lock`, `results-verification-protocol`) remain active  
541 once triggered and explicitly forbid common integrity violations (e.g., post hoc metric changes without au-  
542 thorization; claiming results without fresh verification evidence). Persistent state and locked contracts are ex-  
543 ternalized as on-disk artifacts—not hidden model context—using templates such as `docs/01_intake/research-`  
544 `anchor.yaml` and `docs/03_plan/evaluation-protocol.yaml` (locked after G2), along with experiment logs and  
545 claim–evidence alignment tables. This separation between (i) runtime context (protocol text) and (ii) ver-  
546 ifiable artifacts (contracts and evidence) is what prevents the system from degenerating into prompt-only  
547 compliance and makes the workflow auditable.

## 548 **Procedural workflow specification**

549 The procedural workflow encodes seven phase-gated phases with mandatory temporal ordering and explicit  
550 return paths for backtracking:

551 **Phase 0 (Domain Anchoring)** identifies the research domain, subdomain, research type (Method, Dis-  
552 covery, Tool, or Hybrid), and available resources. Ambiguous inputs trigger clarification rather than assump-  
553 tion.

554 **Phase 1 (Direction Exploration)** executes autonomous literature search (15–30 papers with full-text  
555 retrieval), applies six structured ideation strategies (contradiction mining, assumption challenging, cross-  
556 domain transfer, limitation-to-opportunity conversion, counterfactual reasoning, trend extrapolation), gener-  
557 ates candidate research directions, and refines them through multi-agent brainstorming.

558 **Phase 2 (Problem Validation)** subjects the selected direction to adversarial questioning through a three-  
559 agent deliberation panel, applying a novelty litmus test and feasibility verification. This phase is mandatory

560 for all research types.

561 **Phase 3 (Method/Framework Design)** branches by type. Method projects lock an evaluation protocol  
562 (metrics, datasets, seeds, statistical tests, baselines). Discovery projects design an analysis storyboard with  
563 main and supporting lines, sufficiency criteria, and pre-identified alternative explanations.

564 **Phase 4 (Experiment Execution)** operates in two stages: exploratory validation followed by full ex-  
565 ecution. Method projects require minimum three diagnose–hypothesize–fix–measure iteration cycles. All  
566 experiments use isolated computational environments, fixed random seeds, and logged conditions.

567 **Phase 5 (Results Integration)** compiles outputs, constructs a claim-evidence alignment table, and de-  
568 ploys a three-agent panel for narrative design. Fatal vulnerabilities block progression.

569 **Phase 6 (Paper Writing)** produces modular manuscripts with per-section multi-agent polishing, au-  
570 tonomous reference verification, and full-paper review.

### 571 **Integrity discipline specification**

572 Seven persistent constraints activate at defined trigger points and remain active until project completion:

573 *Metric immutability:* after the evaluation protocol is locked, any modification requires justification, proof  
574 of necessity, and human authorization.

575 *Complete reporting:* all experimental seeds reported with mean and standard deviation; all negative  
576 results recorded; baselines given equal resources; all specified datasets tested.

577 *Claim-evidence alignment:* mapping table linking every paper assertion to supporting evidence; un-  
578 mapped claims flagged for removal.

579 *Alternative-hypothesis exclusion:* systematic confounder assessment before any causal claim.

580 *Reproducibility:* hypothesize–baseline–experiment–verify–interpret cycle with environment logging.

581 *Verification:* fresh computational evidence required before any status claim.

582 *Figure standards:* publication-specific visual quality enforced through pre-defined style profiles.

### 583 **Multi-agent deliberation protocol**

584 Deliberation sessions follow a standardized convergence protocol. A shared scoring rubric is established  
585 before deliberation. Agents assess artifacts independently, revealing assessments simultaneously to prevent  
586 anchoring. After each round, the system checks convergence (all agents PASS, no fatal issues). If not con-  
587 verged, modifications are applied and all agents re-assess the complete artifact (not just changes). Delibera-

588 tion terminates upon convergence, upon reaching the round limit (maximum five), or upon non-convergence  
589 (disagreements escalated to the human researcher).

## 590 **Validation design**

591 We evaluated the formalized methodology through two complementary approaches. First, we applied the full  
592 constraint set to five protocol-constrained research projects spanning diverse scientific domains and research  
593 types. For each project, we documented: (1) domain and research type classification; (2) all phases and gates  
594 completed; (3) integrity constraint compliance (metric stability, seed reporting, negative result recording);  
595 and (4) output artifacts suitable for audit (paper structure, figure inclusion, reference list, and claim–evidence  
596 mapping).

597 Second, to assess whether the constraints themselves (rather than LLM capability alone) account for the  
598 observed rigour, we conducted a direct controlled comparison matched on model, task, tool access, and bud-  
599 get, with the sole difference being activation of the methodology constraints (protocol-free baseline versus  
600 Cursor+Amplify; Section “Controlled study”). Concretely, in the controlled-study project (Project 6) the  
601 agent produced two complete manuscripts under matched conditions with and without the protocol (Sup-  
602 plementary Papers 6A and 6B). We evaluated outputs by auditing generated artifacts rather than subjective  
603 scoring (e.g., manuscript modularity, executable scripts, and the presence of explicit protocol deliverables).

604 **Controlled study (Project 6): protocol-free baseline versus Cursor+Amplify protocol.** To isolate the  
605 contribution of the protocol controller from underlying model capability and tool access, we conducted a  
606 head-to-head controlled study on an identical task using the 1000 Genomes Project 2022 high-coverage  
607 dataset (3,202 individuals, 26 populations). Both conditions used the same IDE environment (Cursor), the  
608 same base model (Claude Opus 4.6), and the same local toolchain (bcftools, PLINK2, Python scientific  
609 stack), under the same constraint of no deep learning model training. The protocol-free baseline condition  
610 (Cursor with no Amplify protocol enabled) was explicitly forbidden from using any Amplify materials and  
611 was prompted to autonomously produce original, publishable scientific findings and a complete manuscript.  
612 The protocol condition (Cursor+Amplify) enabled the phase-gated methodology protocol and emphasized  
613 autonomous execution (running analyses) while producing the protocol deliverables as on-disk artifacts.  
614 Human input was restricted to operational continuation and final compilation requests; no methodological  
615 guidance, scientific direction, or writing edits were provided. Outcomes were evaluated by auditing gener-

616 ated artifacts rather than subjective scoring: manuscript organization (single-file versus modular sections),  
617 analysis scripts and directory structure, and the presence of protocol deliverables (intake anchor, literature  
618 review artifacts, planning documents, and integration blueprint).

### 619 **Code and data availability**

620 The Amplify protocol specification (24 skills, templates, and installation/bootstrapping materials) is released  
621 as open-source software under the MIT licence at <https://github.com/EvoClaw/amplify>. All code used for  
622 the analyses in this paper—including project-specific scripts, figures, and manuscripts for Projects 1–6 (in-  
623 cluding the controlled-study pair, Papers 6A/6B)—is available in the public repositories under the EvoClaw  
624 organization: <https://github.com/orgs/EvoClaw/repositories>.

625 This work introduces no proprietary datasets. Empirical genomics analyses use publicly available 1000  
626 Genomes Project data; simulation-based projects generate data programmatically and include the corre-  
627 sponding code for reproduction in the repositories above.

## 628 **Author Contributions**

629 Y.Z. conceived the study, designed the protocol and controlled study, implemented the project, performed all  
630 experiments and analyses, and wrote the manuscript. J.Z. provided suggestions on biological data analysis  
631 and workflow design, and reviewed the manuscript. All authors approved the final version.

## 632 **Acknowledgements**

633 Amplify and the accompanying project artifacts are released under the MIT licence to encourage reuse,  
634 reproduction, and extension by the research community.

## 635 **References**

- 636 [1] Brian A Nosek, Charles R Ebersole, Alexander C DeHaven, and David T Mellor. The preregistration  
637 revolution. *Proceedings of the National Academy of Sciences*, 115(11):2600–2606, 2018.
- 638 [2] Christopher D Chambers. Registered reports: a new publishing initiative at cortex. *Cortex*, 49(3):  
639 609–610, 2013.
- 640 [3] Kenneth F Schulz, Douglas G Altman, and David Moher. Consort 2010 statement: updated guidelines  
641 for reporting parallel group randomised trials. *BMC medicine*, 8(1):18, 2010.
- 642 [4] Matthew J Page, Joanne E McKenzie, Patrick M Bossuyt, Isabelle Boutron, Tammy C Hoffmann,  
643 Cynthia D Mulrow, Larissa Shamseer, Jennifer M Tetzlaff, Elie A Akl, Sue E Brennan, et al. The  
644 prisma 2020 statement: an updated guideline for reporting systematic reviews. *bmj*, 372, 2021.
- 645 [5] Michael Polanyi. The tacit dimension. In *Knowledge in organisations*, pages 135–146. Routledge,  
646 2009.
- 647 [6] Harry Collins. *Tacit and explicit knowledge*. University of Chicago press, 2019.
- 648 [7] Bruno Latour, Jonas Salk, and Steve Woolgar. *Laboratory life: The construction of scientific facts*.  
649 Princeton university press, 2013.

- 650 [8] Arun James Thirunavukarasu, Darren Shu Jeng Ting, Kabilan Elangovan, Laura Gutierrez, Ting Fang  
651 Tan, and Daniel Shu Wei Ting. Large language models in medicine. *Nature medicine*, 29(8):1930–  
652 1940, 2023.
- 653 [9] Microsoft Research AI4Science and Microsoft Azure Quantum. The impact of large language models  
654 on scientific discovery: a preliminary study using gpt-4, 2023.
- 655 [10] John Jumper, Richard Evans, Alexander Pritzel, Tim Green, Michael Figurnov, Olaf Ronneberger,  
656 Kathryn Tunyasuvunakool, Russ Bates, Augustin Židek, Anna Potapenko, et al. Highly accurate protein  
657 structure prediction with alphafold. *nature*, 596(7873):583–589, 2021.
- 658 [11] Josh Abramson, Jonas Adler, Jack Dunger, Richard Evans, Tim Green, Alexander Pritzel, Olaf Ron-  
659 neberger, Lindsay Willmore, Andrew J Ballard, Joshua Bambrick, et al. Accurate structure prediction  
660 of biomolecular interactions with alphafold 3. *Nature*, 630(8016):493–500, 2024.
- 661 [12] Nathan J Szymanski, Bernardus Rendy, Yuxing Fei, Rishi E Kumar, Tanjin He, David Milsted,  
662 Matthew J McDermott, Max Gallant, Ekin Dogus Cubuk, Amil Merchant, et al. An autonomous  
663 laboratory for the accelerated synthesis of novel materials. *Nature*, 624(7990):86–91, 2023.
- 664 [13] Daniil A Boiko, Robert MacKnight, Ben Kline, and Gabe Gomes. Autonomous chemical research with  
665 large language models. *Nature*, 624(7992):570–578, 2023.
- 666 [14] Kexin Huang, Serena Zhang, Hanchen Wang, Yuanhao Qu, Yingzhou Lu, Yusuf Roohani, Ryan Li, Lin  
667 Qiu, Gavin Li, Junze Zhang, et al. Biomni: A general-purpose biomedical ai agent. *bioRxiv*, 2025.
- 668 [15] Chris Lu, Cong Lu, Robert Tjarko Lange, Jakob Foerster, Jeff Clune, and David Ha. The ai scientist:  
669 Towards fully automated open-ended scientific discovery. *arXiv preprint arXiv:2408.06292*, 2024.
- 670 [16] Yutaro Yamada, Robert Tjarko Lange, Cong Lu, Shengran Hu, Chris Lu, Jakob Foerster, Jeff Clune,  
671 and David Ha. The ai scientist-v2: Workshop-level automated scientific discovery via agentic tree  
672 search. *arXiv preprint arXiv:2504.08066*, 2025.
- 673 [17] Ludovico Mitchener, Angela Yiu, Benjamin Chang, Mathieu Bourdenx, Tyler Nadolski, Arvis Sulovari,  
674 Eric C Landsness, Daniel L Barabasi, Siddharth Narayanan, Nicky Evans, et al. Kosmos: An ai scientist  
675 for autonomous discovery. *arXiv preprint arXiv:2511.02824*, 2025.

- 676 [18] Ali Essam Ghareeb, Benjamin Chang, Ludovico Mitchener, Angela Yiu, Caralyn J Szostkiewicz, Jon M  
677 Laurent, Muhammed T Razzak, Andrew D White, Michaela M Hinks, and Samuel G Rodriques. Robin:  
678 A multi-agent system for automating scientific discovery. *arXiv preprint arXiv:2505.13400*, 2025.
- 679 [19] Juraj Gottweis, Wei-Hung Weng, Alexander Daryin, Tao Tu, Anil Palepu, Petar Sirkovic, Artiom  
680 Myaskovsky, Felix Weissenberger, Keran Rong, Ryutaro Tanno, et al. Towards an ai co-scientist.  
681 *arXiv preprint arXiv:2502.18864*, 2025.
- 682 [20] Juexiao Zhou, Bin Zhang, Guowei Li, Xiuying Chen, Haoyang Li, Xiaopeng Xu, Siyuan Chen, Wenjia  
683 He, Chencheng Xu, Liwei Liu, et al. An ai agent for fully automated multi-omic analyses. *Advanced  
684 Science*, 11(44):2407094, 2024.
- 685 [21] Houcheng Su, Weicai Long, and Yanlin Zhang. Biomaster: Multi-agent system for automated bioin-  
686 formatics analysis workflow. *bioRxiv*, pages 2025–01, 2025.
- 687 [22] OpenAI. Introducing deep research. <https://openai.com/index/introducing-deep-research>, 2025. Ac-  
688 cessed February 2026.
- 689 [23] Hussam Alkaiissi and Samy I McFarlane. Artificial hallucinations in chatgpt: implications in scientific  
690 writing. *Cureus*, 15(2), 2023.
- 691 [24] Dmitry Kobak, Rita González-Márquez, Emőke-Ágnes Horvát, and Jan Lause. Delving into llm-  
692 assisted writing in biomedical publications through excess vocabulary, 2025.
- 693 [25] Matei Zaharia, Andrew Chen, Aaron Davidson, Ali Ghodsi, Sue Ann Hong, Andy Konwinski, Sid-  
694 dharth Murching, Tomas Nykodym, Paul Ogilvie, Mani Parkhe, et al. Accelerating the machine learn-  
695 ing lifecycle with mlflow. *IEEE Data Eng. Bull.*, 41(4):39–45, 2018.
- 696 [26] Mark D Wilkinson, Michel Dumontier, IJsbrand Jan Aalbersberg, Gabrielle Appleton, Myles Axton,  
697 Arie Baak, Niklas Blomberg, Jan-Willem Boiten, Luiz Bonino da Silva Santos, Philip E Bourne, et al.  
698 The fair guiding principles for scientific data management and stewardship. *Scientific data*, 3(1):1–9,  
699 2016.
- 700 [27] Sirui Hong, Mingchen Zhuge, Jonathan Chen, Xiawu Zheng, Yuheng Cheng, Jinlin Wang, Ceyao  
701 Zhang, Zili Wang, Steven Ka Shing Yau, Zijuan Lin, Liyang Zhou, Chenyu Ran, Lingfeng Xiao,

- 702 Chenglin Wu, and Jürgen Schmidhuber. MetaGPT: Meta programming for a multi-agent collabora-  
703 tive framework. In *The Twelfth International Conference on Learning Representations*, 2024. URL  
704 <https://openreview.net/forum?id=VtmBAGCN7o>.
- 705 [28] Qingyun Wu, Gagan Bansal, Jieyu Zhang, Yiran Wu, Beibin Li, Erkang Zhu, Li Jiang, Xiaoyun Zhang,  
706 Shaokun Zhang, Jiale Liu, et al. Autogen: Enabling next-gen llm applications via multi-agent conver-  
707 sations. 2024.
- 708 [29] Marcus R Munafò, Brian A Nosek, Dorothy VM Bishop, Katherine S Button, Christopher D Cham-  
709 bers, Nathalie Percie du Sert, Uri Simonsohn, Eric-Jan Wagenmakers, Jennifer J Ware, and John PA  
710 Ioannidis. A manifesto for reproducible science. *Nature human behaviour*, 1(1):0021, 2017.
- 711 [30] Brian A Nosek, George Alter, George C Banks, Denny Borsboom, Sara D Bowman, Steven J Breckler,  
712 Stuart Buck, Christopher D Chambers, Gilbert Chin, Garret Christensen, et al. Promoting an open  
713 research culture. *Science*, 348(6242):1422–1425, 2015.
- 714 [31] Yilun Du, Shuang Li, Antonio Torralba, Joshua B Tenenbaum, and Igor Mordatch. Improving factuality  
715 and reasoning in language models through multiagent debate. In *Forty-first international conference  
716 on machine learning*, 2024.
- 717 [32] Tian Liang, Zhiwei He, Wenxiang Jiao, Xing Wang, Yan Wang, Rui Wang, Yujiu Yang, Shuming  
718 Shi, and Zhaopeng Tu. Encouraging divergent thinking in large language models through multi-agent  
719 debate. pages 17889–17904, 2024.
- 720 [33] Chi-Min Chan, Weize Chen, Yusheng Su, Jianxuan Yu, Wei Xue, Shanghang Zhang, Jie Fu, and  
721 Zhiyuan Liu. Chateval: Towards better llm-based evaluators through multi-agent debate. *arXiv preprint  
722 arXiv:2308.07201*, 2023.
- 723 [34] Richard E Green, Johannes Krause, Adrian W Briggs, Tomislav Maricic, Udo Stenzel, Martin Kircher,  
724 Nick Patterson, Heng Li, Weiwei Zhai, Markus Hsi-Yang Fritz, et al. A draft sequence of the neandertal  
725 genome. *science*, 328(5979):710–722, 2010.
- 726 [35] Michael Dannemann and Janet Kelso. The contribution of neanderthals to phenotypic variation in  
727 modern humans. *The American journal of human genetics*, 101(4):578–589, 2017.

- 728 [36] Andres M. Bran, Sam Cox, Oliver Schilter, Carlo Baldassari, Andrew D White, and Philippe Schwaller.  
729 Augmenting large language models with chemistry tools. *Nature machine intelligence*, 6(5):525–535,  
730 2024.
- 731 [37] Alireza Ghafarollahi and Markus J Buehler. Sciagents: automating scientific discovery through bioin-  
732 spired multi-agent intelligent graph reasoning. *Advanced Materials*, 37(22):2413523, 2025.
- 733 [38] Eser Aygün, Anastasiya Belyaeva, Gheorghe Comanici, Marc Coram, Hao Cui, Jake Garrison, Renee  
734 Johnston Anton Kast, Cory Y McLean, Peter Norgaard, Zahra Shamsi, et al. An ai system to help  
735 scientists write expert-level empirical software. *arXiv preprint arXiv:2509.06503*, 2025.
- 736 [39] Vlastimil Martinek, Andrea Gariboldi, Dimosthenis Tzimotoudis, Mark Galea, Elissavet  
737 Zacharopoulou, Aitor Alberdi Escudero, Edward Blake, David Cechak, Luke Cassar, Alessandro  
738 Balestrucci, et al. Agentomics: An agentic system that autonomously develops novel state-of-the-art  
739 solutions for biomedical machine learning tasks. *bioRxiv*, pages 2026–01, 2026.
- 740 [40] Yuntao Bai, Saurav Kadavath, Sandipan Kundu, Amanda Askell, Jackson Kernion, Andy Jones, Anna  
741 Chen, Anna Goldie, Azalia Mirhoseini, Cameron McKinnon, et al. Constitutional ai: Harmlessness  
742 from ai feedback, 2022.
- 743 [41] John PA Ioannidis. Why most published research findings are false. *PLoS medicine*, 2(8):e124, 2005.