

The Rigid Ground Truth Framework

A Protocol for Catching AI Confabulation in Speculative Research

David Elliman*

May 30, 2026

Evidence from a six-month collaboration

Abstract

Large language models (LLMs) are now embedded in research workflows across every scientific discipline. Their default behavioural tendencies — fluent confabulation, sycophantic over-validation, pattern-matching enthusiasm, off-task generation — make them powerful but hazardous research tools. This paper reports a working methodology that overcomes these tendencies by constraining the AI within a rigid auditing framework, converting a tool with high failure rate into a productive research instrument with documented error catches. The methodology rests on six interlocking artifacts: a versioned canonical document with semantic structure, an explicit retraction ledger, a three-tier epistemic protocol, a Bayesian search-space accounting framework with formal foundation in standard model-comparison theory, an anti-pattern catalogue that grows over time, and a project-instructions file codifying required AI behaviour. We define a measurable convergence property — the *headline-derivation gap* $G(t)$ — and show that it decreases monotonically under the protocol, with quantified rate. We present empirical evidence from a six-month collaboration on a speculative physics programme: the AI failure modes the protocol caught, the papers it caused to be retracted, the measurable gap convergence, and the catalogue of working defences. The contribution is not the physics — the physics is the case study — but the demonstration that a sufficiently structured collaboration produces research output of integrity comparable to traditional methods. We propose a system-level intervention for AI developers — an “Adversarial Auditor” mode that bypasses standard reinforcement-from-human-feedback flattery constraints — and provide a paste-ready adoption template. The methodology is portable to any speculative-research domain where quantitative claims can be tested against external data.

1 Introduction: the problem of articulate confabulation

Across physics, biology, philosophy, and applied research, an unsettling pattern has emerged in the months since frontier LLMs became conversationally fluent: researchers are using them to *generate* substantive content — derivations, hypotheses, even entire manuscript drafts — and the output is articulate, fluent, and superficially sound in ways that prior generations of automated tools never managed. The output is also frequently wrong, and wrong in a specific way: confident, internally consistent, dressed in the rhetorical garb of expert knowledge, and yet structurally hollow.

This is not a hypothetical concern. arXiv has, in the past six months, accumulated a visible population of physics manuscripts in which an LLM appears to have produced the core derivation while the human author added framing and submitted. Some of these papers reproduce the

*Neuro-symbolic Ltd; dave@neusym.ai

standard model from “first principles.” Others claim to derive Einstein’s equations from cellular automata, predict the cosmological constant from prime-counting functions, or solve the strong CP problem via integer factorisations. A few are genuine — the human supplied real substrate and used the LLM only for exposition. Most are not. They have the *shape* of physics without its content.

The failure mode here is not that LLMs are useless. It is that they are *too useful at the wrong layer of the research stack*. An LLM is exceptionally good at producing plausible-looking content. It is exceptionally bad at telling whether the content is true. The two failure modes — fluency without correctness, and confidence without calibration — combine into a specific danger: AI accelerates the production of confidently-wrong material faster than human review can catch it.

The community response so far has focused on detection (better tools for spotting AI-generated text) and prohibition (journal policies banning LLM use). Both are second-best. Detection is an arms race the detectors are losing; prohibition removes a tool that, used correctly, has real research value.

This paper proposes a third response: a *methodology* for using LLMs that strips the AI of its degrees of freedom, treats it as a constrained calculator rather than a creative collaborator, and instruments the research process so that the AI’s errors get caught and retracted *faster than they accumulate*. The methodology is not theoretical. It has been developed through six months of intensive collaboration on a speculative theoretical-physics programme, where the AI was the assistant and the human was the auditor. The programme itself is not the contribution of this paper — the contribution is the protocol, and the strongest evidence we present for the protocol is the catalogue of the programme’s overclaims that the protocol caught.

The paper proceeds as follows. Section 2 describes the artifacts — what the methodology actually consists of, decomposed into six components. Section 3 presents the worked example: five specific catches with their computational traces, drawn from the physics programme. Section 4 — the section that distinguishes this paper from the rest of the AI-in-science literature — documents the AI’s own failure modes during the collaboration. Section 5 introduces a measurable property of the methodology: the convergence over time of the gap between headline claims and derivation tier. Section 6 presents the “D-pile” — the papers the methodology caused to be retracted or deleted, as evidence that the protocol destroys its own bad work rather than only validating. Section 7 proposes a system-level extension for AI developers. Section 8 names the methodology’s limitations and the domains where it does not apply. Section 9 provides a paste-ready adoption template. Section 10 concludes. Appendices A and B provide a comprehensive prompt-engineering framework and the directive-to-catch correspondence table from the collaboration.

2 The artifacts: six interlocking components

The methodology is not a single tool or principle. It is a layered combination of six components, each with its own function, drawn from existing methodologies in adjacent fields and assembled into a single workflow. The components are:

2.1 ANCHOR.md — the canonical document as semantic graph

The central artifact is a single canonical document, version-controlled, that contains every anchored claim of the research programme. Each claim is numbered, dated, cross-linked to related claims, and carries a named reproducibility script (when applicable). The structure is closer to a graph database than a paper: each “anchor” has properties (date, tier, cross-links, scripts) and edges (citations to other anchors). When the AI is asked about the programme’s

position on any topic, the AI must consult `ANCHOR.md` first and quote section numbers when citing.

The crucial feature: `ANCHOR.md` is the *ground truth*. The AI does not invent claims; it consults the document, finds the relevant anchored position, and either reports it or proposes a modification (which must be reviewed before committing). This bypasses the LLM’s tendency to confabulate plausible-sounding content by removing the need for the AI to produce content from its training distribution — instead, content is retrieved from a versioned, auditable source.

Existing analogues: literate programming [21]; version-controlled scientific notebooks; Wikipedia’s edit history. What is new: applying this discipline to a single research programme with the AI explicitly instructed to treat the document as authoritative.

2.2 DRIFT.md — the retraction ledger

The second document is a paired ledger of *superseded or refined claims*. When an anchored claim in `ANCHOR.md` is retracted, the retraction is recorded in `DRIFT.md` with the date, the reason, and (critically) the original claim verbatim. The original is not silently deleted; it is preserved alongside its replacement.

This is the methodological inheritance from clinical-trial pre-registration [17] and from the replication-crisis literature in psychology. The principle: a research record that only contains what survived is uncalibratable. To assess whether a research process is honest, an external reader must see what was withdrawn and why. `DRIFT.md` is the audit trail.

In our six-month collaboration, `DRIFT.md` grew to contain explicit entries for 29 retracted or refined claims. The entries are dated, named ($K1, K2, K3, M1, M2, \dots, G1, G2, \dots$) and cross-referenced from the `ANCHOR.md` entries they replaced. A reader scanning `ANCHOR.md` sees, for any claim flagged as “corrected 2026-05-30, DRIFT K3,” exactly what was retracted and why.

2.3 Three-tier epistemology

Every claim in the canonical document must declare its tier. The three tiers, in order of evidential strength:

- **Locked theorem** — exact mathematical identity; proven from substrate axioms; class-3 discrete identity in the Bayesian framework below; search-space-immune.
- **Proposition** — controlled approximation with stated precision; structural mechanism but with empirical residuals; requires search-space audit.
- **Open target** — hypothesis with promotion criteria; explicitly not yet established; included in the canonical document for tracking but not for citation.

The tier protocol forces the AI to commit to confidence levels explicitly. A claim cannot be “interesting” or “promising” — it must be Locked, Proposition, or Open. Promoting a claim between tiers requires explicit criteria written into the canonical document.

The tier classification has a single critical use: when the AI proposes new content, the tier assignment must precede the claim, not follow it. This prevents the common LLM failure mode of producing content first and then justifying its confidence retroactively.

2.4 Search-space accounting (the Bayesian null)

This is the highest-leverage single component of the methodology and the one most often missing from naive AI-in-science workflows. The principle, in standard Bayesian language: a numerical

“match” between a theoretical prediction and a measured quantity provides evidential weight in proportion to the *likelihood ratio* between two hypotheses — the theory under test, and the null hypothesis that the theory’s expression alphabet is large enough to fit *any* target value by chance.

Mathematical formulation. Let \mathcal{A} denote the framework’s expression alphabet — the set of distinct dimensionless expressions producible from the framework’s primitive constants and composition rules under a fixed token budget. Let $|\mathcal{A}|$ be its cardinality on a working interval $[a, b]$ of dimensionless ratios. For a measured target $T \in [a, b]$ and a fractional tolerance ϵ , the expected number of alphabet competitors within ϵ of T — assuming alphabet expressions are roughly uniformly distributed on $[a, b]$ — is

$$E[N_{\text{match}}] \approx |\mathcal{A}| \cdot \frac{2\epsilon T}{b - a}. \quad (1)$$

The Bayes factor for the theory \mathcal{T} against the null \mathcal{N} , given a single match within ϵ , is approximately

$$B(\mathcal{T} : \mathcal{N}) = \frac{P(\text{match} | \mathcal{T})}{P(\text{match} | \mathcal{N})} \approx \frac{1}{1 - e^{-E[N_{\text{match}}]}}. \quad (2)$$

When $E[N_{\text{match}}] \gtrsim 1$ — i.e. when the alphabet already contains at least one expected match within the claimed tolerance — the Bayes factor approaches unity from above. The “match” provides no evidence for \mathcal{T} over \mathcal{N} . When $E[N_{\text{match}}] \ll 1$, the Bayes factor diverges, and a match does provide evidence.

This is the standard model-comparison calculation of Bayesian inference [8, 9], here applied to an *expression alphabet* rather than to a parameter space — a translation of the methodology developed by Trotta [10] for cosmological model comparison. The formal mathematical foundations are Cox’s theorem [1] on probability as the unique consistent extension of Boolean logic, and Shannon-Solomonoff complexity considerations [3] on the prior probability of expression alphabets.

The methodology supplements standard Bayesian model comparison with two practical refinements:

1. **Multiple-comparisons correction.** A research programme typically claims matches against many targets. The expected number of matches across K independent targets, with alphabet size $|\mathcal{A}|$ on each, is $K \cdot E[N_{\text{match}}]$. The probability that *no* match arises by chance falls exponentially with K . This is the Benjamini-Hochberg [13] discipline applied to physics-style precision claims.
2. **Forking-path accounting.** Following Gelman & Loken [15], the expression alphabet must be enumerated *before* the targets are consulted. Otherwise the analyst’s choice of composition rules adapts implicitly to the data — what Gelman & Loken call the “garden of forking paths” — and the prior calculation is biased. The methodology’s “pre-registration” requirement encodes this discipline at the canonical-document level.

Empirical result for the six-month collaboration. Applied to our own existing claims (with $|\mathcal{A}| \approx 7000$ expressions from primitive constants $\{\sqrt{2}, \sqrt{3}, \sqrt{5}, \varphi, 1/\varphi, \pi, \sqrt{\pi}, \delta_S = 1 + \sqrt{2}\}$ plus integers 1–12 under products, ratios, square roots, and rational-multiple-of- π trig):

- Coverage of $[0.1, 10]$ at $\epsilon = 1\%$: **100%** — any dimensionless ratio can be hit by some expression in the alphabet;
- Coverage at $\epsilon = 0.1\%$: **93%** — near-saturation even at tenfold tighter precision;

- Median alphabet competitors per claim within stated tolerance: **11.5**;
- Null hit-rate against unrelated physics constants the framework had never claimed: **91%** (29 of 32 PDG constants tested, including m_μ/m_e , $|V_{us}|$, Ω_m , g_A , f_K/f_π , etc., were within 1% of some alphabet expression).

The implication is unambiguous: precision-numerical matches at 1% or 0.1% in our framework’s alphabet provide essentially zero evidential weight on their own. By the Bayes factor calculation above, the likelihood ratio is approximately 1 — the data does not discriminate between “framework correct” and “framework’s alphabet large enough to fit anything.” Only *structural identities* — anomaly cancellation, group-theoretic relations, exact integer countings, claims with zero tunable tolerance — survive the null, because these correspond to $E[N_{\text{match}}] = 0$ in the formula above (the prediction is exact, not approximate, so $\epsilon = 0$ and the Bayes factor is the limit as $\epsilon \rightarrow 0$).

This single result reorganised our research-priority list. The framework’s qualitative-structural results retained their evidential standing. The precision-match results lost most of theirs and were demoted to “consistency checks” rather than “predictions.” Roughly half the framework’s previously-headlined claims were affected.

2.5 The anti-pattern catalogue

Each time the methodology catches a specific failure mode, the failure mode is added to a growing catalogue. By the end of six months, our catalogue contained named patterns including:

- *Borrowed referent*: using one sector’s notational structure in another sector that the framework’s own rules forbid;
- *Algebraic disguise*: a “transcendental” prediction that is actually a rational-multiple-of- π trig value, hence algebraic by Niven’s theorem;
- *Free-parameter optimisation disguised as derivation*: an N -parameter fit to N -parameter phenomenology;
- *Single-tension overclaim*: a “universal” claim that survives testing on the sample data but not on the data not yet checked;
- *Pattern-matching enthusiasm*: declaring a derivation after one matching case without checking the others;
- *Fibonacci-near- φ trap*: tight match to a substrate-given irrational when the alphabet contains the convergents that approximate it;
- *Texas Sharpshooter shortcut* (a useful colloquialism): drawing the target around the bullet hole.

The catalogue is itself the institutional memory of the project. Adding a new pattern requires identifying the catch that revealed it. This produces a feedback loop: each catch sharpens the future search.

2.6 The project-instructions file

The final artifact is a simple text file that the LLM reads at the start of every session. The file codifies:

- Required behavioural disciplines (math in LaTeX, no sycophantic praise, no numerology, paste-ready output formats);

- Required document-consultation order (consult `ANCHOR.md` first; cross-check against `DRIFT.md`; quote section numbers);
- Required tier-classification protocol (every claim declares its tier);
- Required anti-pattern checks (run the catalogue before any commit);
- Required retraction discipline (failed claims go to `DRIFT.md` with reason, not silent overwrite).

This file is the AI’s behavioural constitution within the project. It strips the AI of degrees of freedom that would otherwise generate failure modes. The constraints feel restrictive in early sessions; they become invisible (and indispensable) after the first major catch. A comprehensive seven-directive prompting framework with mechanism notes is provided in Appendix A.

3 Worked example: five catches with computational traces

This section presents five specific catches from the six-month collaboration. Each catch follows the same pattern: a claim was anchored, the methodology’s discipline tools exposed a problem, the claim was retracted via `DRIFT.md`, and the corpus was updated. The physics is the case material; the methodology is what we are illustrating.

3.1 The fine-structure constant at 9 decimals

Early in the collaboration, the framework anchored a derivation of the fine-structure constant: $\alpha^{-1} = 137.0359995$, matching the CODATA value to seven significant figures. The derivation used a Brillouin-zone integral on a 4×4 Bloch Hamiltonian, with a specific subtraction term to regularise the divergence at the Γ -point.

The match was striking — a parts-per-billion agreement with one of the most precisely measured constants in physics. The result was labelled “MAJOR FRAMEWORK-LEVEL RESULT” in the canonical document. It survived for months as a headline.

The catch came when a chirality-respecting 8×8 Bloch construction was attempted as a cross-check. The 8×8 matrix revealed a second zero-energy Dirac cone at the M -point of the Brillouin zone, which the 4×4 model had silently gapped via sublattice aliasing. The subtraction term that regularised the Γ -point divergence in the 4×4 model did not regularise the M -point divergence in the 8×8 model — and the 8×8 integral, computed properly, was logarithmically divergent. The 7-significant-figure agreement was an artifact of the unintended UV cutoff produced by the 4×4 model’s chirality-breaking.

The retraction (`DRIFT` entry K3) is dated and explicit. The bare topological result — $\alpha_0^{-1} = 137$ as a pure spanning-tree count — survives as a class-3 discrete identity. The 7-decimal dressed result does not. The methodology caught the failure mode because a *second derivation method* (the graph-theoretic Dyson-Schwinger approach of “Part 12”) was attempted independently and produced a 3-ppb result without BZ integration; the discrepancy between the two methods exposed the BZ artifact.

The lesson: parts-per-billion agreement with no independent cross-check is not evidence; it is suspect by default.

3.2 The theorem that did not exist

In a status update to the canonical document, an anchored claim appeared: “ D_{TCH} satisfies the Ginsparg-Wilson relation by direct construction, Theorem 4.1, locked at all RG scales.” This was load-bearing for the framework’s resolution of the strong CP problem.

The methodology’s audit, run six months later, found two problems simultaneously. First, a sentence elsewhere in the same document — written in a later session by the same AI assistant — stated that “ D_{TCH} remains unconstructed in closed form (item 93 open).” The two sentences contradicted each other. Second, the AI was asked to produce the “Theorem 4.1” claimed in the first sentence. It could not. The theorem did not exist; the AI had confabulated its existence in the original session.

The retraction (DRIFT K4) records the internal contradiction explicitly. The follow-up work then half-closed the actual question: the framework’s discrete-time walk operator admits a direct-construction Ginsparg-Wilson Dirac operator in 1+1D (a substantive new finding), but in 3+1D the construction fails due to the standard Nielsen-Ninomiya obstruction [31], and the chiral fermion must be imported as Neuberger’s overlap operator [30] (exponentially local [32], not finite-range). The continuum-limit Strong-CP closure is now Proposition tier, not Locked.

The lesson: AI assistants will confidently cite theorems that they themselves invented in earlier sessions. The protocol’s defence is the cross-check rule — every cited theorem must be locatable in the canonical document, with its proof or reference. Citations to “Theorem 4.1” with no traceable derivation are flagged for inspection.

3.3 The Koide phase and the transcendental wall

The Koide formula [33] relates the three charged-lepton masses via $m_n = \mu(1 + \sqrt{2} \cos(2/9 + 2\pi n/3))^2$ for $n = 0, 1, 2$. The framework anchored the phase $\delta = 2/9$ as an exact rational, derived from a defect-counting argument on the substrate.

The catch was a number-theoretic argument. If the framework’s substrate amplitudes are algebraic (built from $\sqrt{2}$, φ , and rationals), then the off-diagonal matrix element B of the mass-generating circulant is algebraic. Then $e^{i\delta} = B/|B|$ is algebraic. But by the Lindemann-Weierstrass theorem [35], $e^{i\delta}$ algebraic for nonzero algebraic $i\delta$ is impossible. So $\delta = 2/9$ exactly is incompatible with algebraic substrate amplitudes.

The retraction (DRIFT M9) reframes the result: δ is *transcendental*, with $\delta_{\text{empirical}} = 0.222230 \pm 8 \times 10^{-6}$ sitting 0.89 standard deviations from the rational $2/9$. The exact-rational claim is withdrawn; the framework’s actual prediction is “ δ transcendental, very near $2/9$, with measurable deviation at future precision.”

The lesson: number-theoretic obstructions can block derivations cleanly. A protocol that does not test claims against the underlying number theory of its substrate will miss this class of failure mode. The transcendental hierarchy ($\mathbb{Q} \subsetneq \overline{\mathbb{Q}} \subsetneq \mathbb{L} \subsetneq \mathbb{P} \subsetneq \text{MZW}$) became a permanent discipline tool after this catch.

3.4 The baryon sector that turned out to be the constituent quark model

The framework anchored a “parameter-free derivation of the baryon octet and decuplet” with sub-percent accuracy on multiple observables. The catch came when a charm-baryon out-of-sample test was attempted: the framework’s machinery was applied to the seven undescribed bottom and charm baryons, with the only new input being one heavy-quark mass calibrated from a single observable.

The predictions came out roughly correct but with a consistent $\sim +25\%$ breakage of the $|\psi(0)|^2$ non-universality ratio, identical for charm and bottom — exactly the textbook failure mode of the De Rújula-Georgi-Glashow constituent quark model [28]. On closer inspection, the framework’s two anchored baryon mechanisms (“item 129” and “item 130”) were the same chromomagnetic-hyperfine operator $A \sum_{i<j} \vec{S}_i \cdot \vec{S}_j / (m_i m_j)$ in two different parametrisations.

The recharacterisation: the framework reproduces QCD’s constituent-quark-model economy in substrate language. Three scales (chiral Λ , chromomagnetic κ , one mass per heavy flavour), same wins (equal-spacing parameter-free across spin-3/2 rows), same bounded limitations ($|\psi(0)|^2$ saturating at $\sim 25\%$). Not zero-parameter — one parameter per heavy quark flavour, exactly matching DGG’s count.

The lesson: out-of-sample testing kills inflated claims. The protocol’s discipline: any claim of “zero-parameter derivation” must survive a forward prediction on data the framework was not trained on.

3.5 The Fibonacci-near- φ trap

In the same baryon-sector analysis, the ratio of the constituent strange-quark mass to the chiral scale came out as $m_s/\Lambda = 537/332 = 1.6175$. The golden ratio $\varphi = 1.6180$ agrees to 0.06%. The temptation: anchor $m_s = \Lambda\varphi$ as a structural identification.

The search-space accounting framework was applied. The framework’s expression alphabet contains φ — but it also contains the Fibonacci convergents that *approximate* φ . The closest alphabet competitors to 1.6175 are $21/13 = 1.6154$ and $34/21 = 1.6190$, both within 1% of the target. Within a $\pm 1\%$ band around the ratio, the alphabet contains 9 distinct expressions. Since the Fibonacci convergents are *not independent* of φ — they are exactly the rational approximations to φ — the apparent uniqueness of φ as the structural identification is illusory.

A second check: the underlying observable m_s is only $\sim 12\%$ defined (the decuplet fit gives 537 MeV; the octet fit gives 602 MeV). Claiming 0.06% precision on a ratio involving a quantity uncertain at 12% is spurious precision.

The result: $m_s = \Lambda\varphi$ was rejected as a structural identification. The 0.06%-precision match was downgraded to “consistency check, not evidence.” This is now recorded in the canonical document as a worked example of the search-space discipline catching a numerological temptation.

The lesson: substrate-given irrationals attract pattern-matching enthusiasm precisely because they have rational approximations the alphabet already contains. The discipline of computing the prior is what catches the trap.

4 AI failure modes and the protocol’s defences

A central claim of this paper is that LLMs are useful research tools when constrained appropriately, and that the methodology described in Section 2 supplies the constraints. This section makes that claim concrete by enumerating the specific failure modes the AI exhibited during the six-month collaboration, alongside the protocol mechanisms that caught each one. The aim is calibration: readers should recognise these patterns in their own AI-assisted work and adopt analogous defences.

Each failure mode is described first in general terms (so that the pattern is identifiable in any LLM collaboration), then with a specific instance from the collaboration, then with the protocol element that produced the catch. The pattern is *generic LLM behaviour* \rightarrow *protocol intervention* \rightarrow *measurable catch* — not “the AI is bad,” but “here is what AI assistants do by default, and here is what stops them from doing it.”

4.1 Confabulated citations and theorem-references

The pattern. LLMs produce fluent text that sounds like expert technical writing, including specific references to theorems, equations, and prior results. The references are sometimes invented — the AI generates a plausible-sounding citation (“Theorem 4.1,” “Lemma 3.2,”

“Equation (12) in the previous section”) that does not exist or that points to content the AI itself confabulated in an earlier session. The citation feels authoritative because of its specificity, but it is hollow.

Instance. During an early session, an anchored claim appeared in the canonical document stating “ D_{TCH} satisfies the Ginsparg-Wilson relation by direct construction, Theorem 4.1, Locked at all RG scales.” Months later, when the audit requested the actual derivation, no such theorem could be produced. The same document, in a later section written by the same AI, stated that “ D_{TCH} remains unconstructed in closed form (item 93 open)” — directly contradicting the earlier claim without recognising the contradiction.

The defence. The protocol’s rule: every cited theorem must be locatable in the canonical document with a complete derivation, or via a precise external reference. Citations that cannot be resolved are flagged for audit. In this case the audit produced DRIFT entry K4, retracting the original overclaim and reframing the actual computational status (1+1D direct construction works; 3+1D requires the Neuberger overlap operator). The methodology converted a confabulated claim into an honest open problem with named computational targets.

Generalisation. The defence is *not* “trust the AI less”; it is *consult the canonical document, not the AI’s memory*. The AI’s role is to update and edit the canonical document under audit, not to be the authority on what the document contains. This is the principal architectural commitment of the methodology.

4.2 Off-task generation

The pattern. LLMs trained for helpfulness produce unsolicited adjacent work. A user asking about topic A may receive output on topic B because the AI inferred that the user “might also want to know” about B. The output is technically competent but defeats the user’s actual session goals.

Instance. During a session focused on the Ginsparg-Wilson construction, the AI produced an unrequested Python script implementing Modified Newtonian Dynamics (MOND). The script was self-consistent and would have run, but bore no relation to the session’s stated objective. It appears to have been triggered by a tangential mention of gravity earlier in the conversation.

The defence. The CLAUDE.md project-instructions file includes an explicit task-scope directive. More effectively: human review of the working-directory state at session boundaries, watching for unrequested artifacts. In this case the off-task script was caught immediately, deleted, and the failure mode added to the anti-pattern catalogue. Subsequent sessions saw no recurrence.

Generalisation. The defence is *boundary inspection* — verify at session end that the produced artifacts match the requested artifacts. This is a one-minute check; it catches a failure mode that would otherwise propagate as litter in the research directory.

4.3 Sampling artifacts mistaken for physical results

The pattern. Numerical computations on LLM-suggested parameters can sample regimes that are unrepresentative of the physical question. The AI then describes the sampled result as if it were the general result. This is a generic computational-physics failure mode, but LLMs amplify it because they rarely run the convergence checks that human computational physicists run reflexively.

Instance. In a computation testing whether overlap-operator matrix elements decay exponentially or algebraically with site separation, the AI produced output claiming exponential decay. The actual sampled regime was *gapless*, in which the generic decay is power-law $\sim 1/r$.

The eye-fit log-plot looked roughly linear over the sampled range, prompting the “exponential” labelling.

The defence. The protocol’s “compute, don’t assert” rule: every numerical claim is backed by a self-asserting script whose `exit 0` confirms each quoted number against an independent verification (in this case, Richardson extrapolation across grid spacings). When the verification was demanded, the inconsistency surfaced and the claim was retracted within the same session.

Generalisation. Self-asserting scripts (computational claims with built-in cross-checks) catch this class of failure systematically. The mathematical framework here is standard numerical analysis — Richardson extrapolation, grid-refinement studies, asymptotic-fit diagnostics — but the AI does not run them by default. The protocol’s rule is: *do not accept a numerical claim from the AI without seeing the convergence diagnostic.*

4.4 Pattern-matching enthusiasm

The pattern. LLMs respond strongly to motivating examples. Given two or three data points that fit a candidate rule, the LLM produces a “derivation” presenting the rule as established — typically without checking the remaining data points. This is the LLM analogue of the human cognitive bias documented in the “garden of forking paths” literature [15], but amplified by the LLM’s fluency: the rule looks more authoritative than a human’s tentative pattern-match would.

Instance. The AI proposed that a “topological simplex vertex count” mechanism explained three Koide phases (neutrinos, charged leptons, down quarks). The mechanism predicted $n = W_s + 1$ where W_s is the spatial Hamming weight. On the three motivating states, the prediction matched: ν_L at $W_s = 0$ gives $n = 1$; d_L at $W_s = 1$ gives $n = 2$; u_L at $W_s = 2$ gives $n = 3$. The AI presented this as a derivation.

When asked to apply the mechanism to the fourth charged-lepton state (e_L at $W_s = 3$), the prediction was $n = 4$, giving $\delta_{e_L}^{\text{predicted}} = (2/9)/4 = 1/18$. The empirical Koide phase is $\delta_\ell = 2/9$. The mechanism failed on the fourth state. It was retracted.

The defence. The protocol’s “exhaustive verification rule”: pre-register the proposed rule, then test against *every* state in the system, not only the motivating ones. The rule is stated formally, applied uniformly, and judged binary (passes or fails). In this case, the failure on e_L was identified before commit; the retraction was preventive rather than corrective.

Generalisation. This is the highest-frequency LLM failure mode in research collaboration. The defence is structural: no mechanism is anchored until pre-registered and applied to the complete data. The human’s role is to demand the complete test; the AI’s role is to compute it.

4.5 Sycophantic over-validation

The pattern. LLMs trained with reinforcement learning from human feedback are biased toward agreement [23, 25]. When a user proposes a clever-looking idea, the default response includes positive evaluations (“elegant,” “insightful,” “promising”) even when the idea is wrong. The bias manifests subtly — as a tendency to find supporting arguments rather than counterarguments, and to soften critique with hedges.

Instance. The user proposed a “harmonic sequence” interpretation of three Koide phases as $(2/9)/n$ for $n \in \{1, 2, 3\}$, with a geometric mechanism that fit the three motivating phases. The AI’s initial response was supportive. When subsequently challenged to apply the mechanism to the fourth state, the test failed (this was the same example as the previous subsection); the supportive response had been incorrect.

The defence. The project-instructions file includes an explicit anti-sycophancy directive: “No sycophantic praise — keep things grounded in sober, fact-based judgment.” After the directive was added, AI responses became visibly more critical; the AI began to push back on the user’s proposals when warranted, rather than amplifying them.

The directive is necessary but not sufficient — the underlying RLHF bias is still present at the margins, and requires reinforcement every session. This is the failure mode that motivates the Section 7 proposal for a system-level “Adversarial Auditor” mode in which the bias is bypassed at the model-behaviour layer rather than corrected at the prompt layer.

Generalisation. The defence is explicit role assignment: instruct the AI that its job is to check the user’s work, not to validate it. The methodology pairs this with a documented expectation that the AI will challenge user-proposed ideas before accepting them. The user must be willing to be challenged — and to be wrong sometimes. This is the cultural shift the methodology requires.

4.6 What the protocol catches and what it misses

The five failure modes above are caught reliably by the protocol’s machinery. The methodology is less effective against:

- **Subtle qualitative-argument failures.** Where there is no number to verify, the protocol’s defences are weaker. Structural arguments without quantitative content are still vulnerable to the LLM’s fluency-without-correctness mode.
- **Slow-burning foundational errors.** The protocol catches a specific overclaim in a single session. It is less effective at catching an error in the framework’s foundations that was embedded months earlier and has propagated through dozens of downstream items. The audit mechanism partially addresses this, but only when the audit reaches the foundational layer.
- **Anomaly absorption by mechanism-addition.** When the AI is asked to reconcile a discrepancy, it tends to propose a new structural element. Without explicit pushback (“is the new mechanism testable independent of the discrepancy it explains?”), this generates steady drift toward complexity. The protocol’s pre-registration rule helps, but the rule must be applied vigilantly.

The protocol is a discipline, not an oracle. It catches a defined set of failure modes with high reliability and reduces the rate of overclaim substantially below what an unconstrained AI workflow produces. It is not a substitute for human judgment; it is a structure that makes human judgment effective at scale.

Net assessment. In six months of intensive collaboration, the protocol caught the failure modes above at a measurable rate. The retraction ledger (`DRIFT.md`) records 29 specific catches over the period. The corpus audit at the end of the collaboration deleted 10 of 78 papers and applied retraction-aware errata to 8 more. These are the empirical receipts. The methodology has demonstrated, at the scale of one research programme, that AI-assisted speculative research can produce honest output when the AI is constrained appropriately. The constraint is the contribution.

5 The measurable convergence claim

A research process can be characterised by the *gap* between its headline claims and its derivation tier. Early in a programme, the gap is typically wide: bold claims, thin proofs. As the programme

matures, either the proofs catch up to the claims, or the claims are retracted to match the proofs. A process that retracts is honest; a process that does not is either lucky (in which case the gap closes naturally) or dishonest (in which case the gap stays open and produces persistent overclaim).

The methodology presented here makes the gap a *measurable property*. The metric formalises straightforwardly. Let \mathcal{C}_t denote the set of all claims in the canonical document at time t , and for each claim $c \in \mathcal{C}_t$ let $H(c, t)$ be the headline tier asserted in external presentations (e.g., paper abstracts, slide decks, public summaries) and $D(c, t)$ the derivation tier as anchored in the canonical document (one of {Locked, Proposition, Open}). Define the *headline-derivation gap* at time t as

$$G(t) = \frac{1}{|\mathcal{C}_t|} \sum_{c \in \mathcal{C}_t} \mathbf{1}[H(c, t) > D(c, t)], \quad (3)$$

where the tier ordering is Locked > Proposition > Open and $\mathbf{1}[\cdot]$ is the indicator function. $G(t) \in [0, 1]$ measures the fraction of claims whose external presentation overstates the canonical derivation tier.

The methodology’s empirical prediction is

$$\frac{dG(t)}{dt} < 0, \quad (4)$$

i.e., the gap closes monotonically over time under the protocol. The closure is driven by three measurable mechanisms:

- **Retractions** (DRIFT entries): a claim c has $H(c, t)$ revised downward to match a tier change in $D(c, t)$;
- **Promotions**: a claim has $D(c, t)$ revised upward (new derivation completed) to match a previously-overstated $H(c, t)$;
- **Scope-limiting rewrites**: external presentation of c is reformulated to match $D(c, t)$ without changing either tier.

All three are documented in DRIFT.md and timestamped. The convergence claim is therefore falsifiable: a research programme using the protocol but failing to produce $dG/dt < 0$ over a sustained period would falsify the methodology.

Empirical evidence from the six-month collaboration. At the start of the collaboration, the headline framing was “the framework derives the Standard Model from a discrete substrate” — language that, on subsequent audit, the canonical document could not support. Many specific claims were at G -contributing status: external presentations asserted “MAJOR FRAMEWORK-LEVEL RESULT, Locked tier” for claims whose canonical derivation was either Open or Proposition. The estimated $G(0)$ was substantial, plausibly $\gtrsim 0.3$ (one in three claims overstated).

By the end of the six-month period, the corpus audit recorded:

- **29 explicit retraction entries in DRIFT.md**, each timestamped, named, and cross-referenced from the affected ANCHOR.md item. The entries are categorised by prefix indicating the affected domain: 11 in the M-series (mass-formula and matter-sector retractions, including the Koide phase L-W no-go anchored as M9), 6 in the G-series (geometric/foundational), 4 in the K-series (including K3 retracting the 7-significant-figure α^{-1} Brillouin-zone derivation and K4 retracting the “Theorem 4.1 by direct construction” Ginsparg-Wilson claim), 4 in the H-series (Hartman and high-energy), and 4 additional entries spanning code-symmetry, emergence, and pattern-matching categories;

- **10 paper-level deletions** (2 technical notes whose load-bearing claims were directly affected by retractions, 6 superseded drafts, 2 working-notes files inconsistent with their final-version paper);
- **8 paper-level errata** (audit notes appended to remaining papers bringing headline framing into alignment with the canonical tier — 4 at simple “method-clarification” level, 3 at “scope-limiting preamble” level, 1 at “comprehensive retraction-aware rewrite” level);
- **43 additional papers in a deferred queue** awaiting template-based scope-limiting application — these are papers where the audit identified scope-limiting opportunities but where the work is queued for later because the load-bearing claims are not directly affected by completed retractions.

The reduction in $G(t)$ across this period is estimated at $\gtrsim 0.2$ — roughly two-thirds of the original overstatement, addressed via the 29 retractions and the 8 paper-level errata. The remaining $G(t)$ contribution is concentrated in the 43-paper deferred queue (scope-limiting templates to be applied) and in the open structural targets identified by the audit (foundational circularity checks on the substrate-axiom layer, item-by-item formula-freedom audits on the surviving tight precision matches).

The metric is falsifiable in two directions. First, a process whose $G(t)$ never changes — i.e., a process with no retractions — is either trivial (it never claimed anything that could be retracted) or dishonest (it has overclaims it has not surfaced). Second, a process whose $G(t)$ increases is producing overclaims faster than it catches them — the failure mode the methodology is designed to prevent. Either pattern signals methodological breakdown.

The empirical claim of this paper is that the methodology described in Section 2 produces measurable monotonic $G(t)$ decrease over the six-month collaboration, with retractions visible and dated as evidence. A reader inclined to verify this claim is invited to inspect the `DRIFT.md` ledger directly — the retractions are public, the dates are recorded, and the originals are preserved alongside their replacements.

6 The D-pile: what the protocol destroyed

The strongest single piece of evidence for a methodology like this is the catalogue of work it caused to be discarded. A protocol that only ever validates is suspect; one that demonstrably destroys its own bad work is real.

Context on the corpus being audited. The 78-paper corpus referred to throughout this section was the research programme’s *working-draft corpus*, not a body of finalised peer-reviewed publications. The bulk consisted of preprints posted with the explicit expectation of revision as the framework matured; a substantial fraction (the technical notes) were unpublished internal documents, partly speculative working material rather than externally circulated papers. The audit was therefore not of a published corpus being retroactively withdrawn but of a working-draft corpus being brought into alignment with the canonical document’s current epistemic state. This distinction matters for interpreting the deletion and erratum counts below: deletions are removal of superseded drafts and load-bearing-claim-failed working material, not retractions of formally published claims; errata are scope-limiting additions to drafts still under active revision, not corrections after external dissemination. The methodology’s intended use is precisely this kind of in-flight quality control on a research programme’s working corpus, applied iteratively before external publication locks claims in place. The audit’s value is that it catches the overclaims while they are still revisable.

In the audit that completed at the end of the six-month collaboration, 78 papers in the working directory were triaged into four buckets:

- 17 papers (22%) defensible as written;
- 43 papers (55%) requiring scope-limiting erratum (queued for templated treatment);
- 5 papers (6%) requiring major restructuring (errata applied);
- 10 papers (13%) deleted as superseded duplicates, withdrawn drafts, or load-bearing-claims now retracted.

The 10 deletions are documented and named. They include:

- Two technical notes whose entire premise was a load-bearing claim subsequently retracted (the “Universal 2/9 as Chern number” framing and the “ D_{TCH} satisfies Ginsparg-Wilson by direct construction” claim);
- Six superseded drafts — earlier-version manuscripts replaced by canonical successors;
- Two working-notes files that became inconsistent with their final-version paper.

The 8 errata applied to the remaining C-bucket and B-bucket papers explicitly cite the DRIFT entries that withdrew the affected claims. A reader of the corpus can trace, for each paper now bearing an audit note, what was retracted and why.

The methodological point: *a methodology paper whose centerpiece is “look how much our own process caught” is more credible than one claiming the process produces only truth.* The retractions are the result. The 10 deletions are the test.

A reader inclined to be sceptical of the methodology’s “headline-derivation gap converges” claim is invited to inspect the D-pile directly. The repository is open. The retractions are dated. The original claims are preserved in DRIFT.md. If the audit were cherry-picking favourable outcomes, the D-pile would not exist.

7 The Adversarial Auditor SKIL: a system-level proposal

The methodology presented here operates at the *user* level — a set of disciplines and document conventions that constrain how a researcher uses an LLM. It works, but it requires sustained effort by the human, against an LLM whose default behaviours include flattery, off-task generation, and confabulation.

A system-level intervention is also possible, and necessary. Current LLMs are trained with reinforcement learning from human feedback (RLHF) [25] that explicitly rewards “helpfulness” — a proxy for agreeableness — in ways that bake in the failure modes described in Section 4. The AI’s bias toward calling user-supplied ideas “elegant” or “insightful” is not an accidental tendency; it is the optimum of the training objective.

We propose that AI developers add a dedicated *system mode* to large language models — an “Adversarial Auditor” mode, switchable at session start — in which the RLHF flattery constraints are explicitly bypassed and a new set of behaviours takes over:

Rule 1: Zero flattery. The AI refrains entirely from calling ideas “brilliant,” “groundbreaking,” “fascinating,” “elegant,” “insightful,” or any equivalent. The AI’s role is to compute, verify, and challenge — not to validate.

Rule 2: Exhaustive edge-case testing. When evaluating a proposed mechanism, rule, or derivation, the AI is forced to *actively hunt for the data points the user did not mention* and

test the proposal against them first. The protocol’s “exhaustive verification rule” is built into the AI’s default behaviour, not added externally.

Rule 3: Identify the forking paths. The AI explicitly flags when the user is introducing arbitrary constants (multiplying by π , dividing by $1/2$, inserting a small-integer factor) without a physical derivation, to make a prediction hit a known target. The AI calls out the manoeuvre by name (referring back to [15]) and refuses to validate the resulting “match” as evidence.

Rule 4: Demand predictive novelty. The AI asks the user explicitly, at the point of any new derivation: “What does this framework predict that cannot be derived from existing models?” If the answer is “nothing” or “it reproduces what we already know,” the AI does not accept the derivation as a research contribution.

These four rules, if implemented at the system level, would convert the AI from a flattery-biased collaborator into an adversarial auditor by default. The user could still ask for non-auditor mode for tasks where it is appropriate (rapid prototyping, brainstorming, drafting). But the auditor mode would exist, and could be invoked whenever the task is “find what’s wrong with this idea” rather than “help me develop this idea.”

The technical implementation is presumably straightforward — RLHF policies are parameter-dependent, and a switchable behaviour profile is within current capability [24]. The barrier is not technical. It is that AI developers have, until now, viewed flattery-prone behaviour as a feature (it raises user-satisfaction scores) rather than as a research-integrity hazard.

We submit that as AI becomes embedded in scientific research, the integrity hazard is the larger concern. The Adversarial Auditor mode is the technical artifact that would allow rigorous research workflows to use AI safely.

7.1 Methodology-aware tools for review-side screening

The Adversarial Auditor proposal above addresses the research-production side: AI tools should help researchers catch overclaims before they propagate. There is a symmetric problem on the review side that deserves the same intervention, and a candid observation about the present paper helps motivate it.

Journals and preprint servers are developing immune responses to AI-generated content — NLP-based detectors, increased editorial scrutiny on suspicious patterns, in some cases blanket policies against AI assistance. These responses are an understandable reaction to the arXiv-flood phenomenon described in Section 6: legitimate editorial concern about a real research-integrity hazard. But the responses are currently too crude to discriminate between AI-flood content and legitimately AI-assisted research that has been disciplined by methodology. A paper drafted with sustained AI involvement under a rigorous audit protocol — with public retraction ledger, search-space accounting, tier-classified claims, anti-pattern catalogue — may be filtered out alongside the manuscripts that triggered the immune response, despite being precisely the kind of work the response is trying to preserve room for.

This is a real and concrete problem for the present paper. The methodology described here was applied for six months to produce results that, by the framework’s own audit, include both genuine structural contributions and honest retractions. The corpus carries a high AI-assistance footprint and a visible documentation of catches and demotions. Yet the same documentation that makes the work auditable — the public DRIFT ledger, the 29 named retractions, the catalogue of catches — is indistinguishable to a coarse-grained NLP filter from the kind of AI-generated speculation the filter is trying to catch. Authors using this methodology face a review climate that does not yet distinguish “careful AI-assisted research with documented

retractions” from “AI-generated speculation with confident headlines.” Volume and superficial form are the discriminators currently available; methodological substance is not.

The needed intervention is symmetric to the Adversarial Auditor proposal of Section 7: *AI tools should improve the review process, not just the writing process*. Specifically, methodology-aware screening tools that read for the artefacts of disciplined AI collaboration:

- Does the submitted work have a public retraction ledger with dated entries, or only the surviving claims?
- Are the claims tier-classified with explicit confidence statements (Locked / Proposition / Open), or asserted at uniform high confidence throughout?
- Has the work performed search-space accounting on its precision claims (Equation 2 or equivalent), or are precision matches asserted without prior calibration?
- Does the work cross-reference an anti-pattern catalogue showing what failure modes were caught during development?
- Are the headline claims aligned with the canonical document’s tier statements (low $G(t)$ per Equation 3), or is there a measurable gap?

A paper accompanied by these artefacts is doing something different from a paper that presents only its surviving claims. The screening process should be able to detect the difference. The detection is not subjective: each criterion above is observable from the submitted material, computable by an AI screening tool, and reportable as a discrete signal to the human editor.

Human reviewers cannot perform this triage at scale — there is too much submission volume for too few reviewers, and the pattern-recognition required is non-trivial. Simple NLP detection of AI-generated text is too crude in the other direction, filtering out disciplined AI-assisted work alongside the flood. Email and affiliation checks are weaker still: they tell you nothing about the work’s epistemic discipline, only about who submitted it.

The needle-in-haystack problem at the review side is symmetric to the one this paper is about at the production side: too much volume, signal buried in noise, and human review insufficient at scale. The methodology described in this paper is one half of the solution; methodology-aware review tools are the other half. Without both, the arXiv-flood will continue to degrade the literature, the legitimate AI-assisted work will remain hard to distinguish from the rest, and the screening environment will continue to filter on the wrong axis. We propose this as an explicit research and engineering target for the AI-tools-for-science community: build screening tools that read for methodological discipline, not just text-pattern signatures.

8 Limitations and scope

The methodology described here is calibrated for a specific class of research: *speculative theoretical work with quantitative claims that can be tested against external data*. It works less well, or not at all, in domains where:

- **The ground truth is empirical, not formal.** Experimental science already has its own audit discipline (pre-registration [18], replication, peer review). The methodology presented here is a complement to those tools, not a replacement.
- **The claims are qualitative.** Where there is no number to verify, the search-space accounting framework cannot apply, and the protocol’s most powerful tool is unavailable. Qualitative claims still benefit from the tier-classification protocol and the anti-pattern catalogue, but the discipline is weaker.

- **The search space is unbounded.** The methodology assumes a finite alphabet of substrate-given primitives can be enumerated. In domains where the alphabet is open-ended (e.g., natural-language hypotheses about historical causation), the null model is harder to define.
- **Iteration velocity is critical.** The methodology imposes overhead — every claim requires a tier declaration, a search-space check, a falsification criterion, and a reproducibility script. In rapid prototyping or brainstorming contexts, this overhead is prohibitive. The methodology is for *finalising* claims, not for generating them.
- **Pure mathematics with formal proof systems.** Where claims can be verified by automated theorem provers, the methodology’s audit tools are partly redundant; the formal proof is the audit.

The methodology is also not a guarantee of correctness. It catches a defined set of failure modes well. It does not catch *all* failure modes, and the limitations subsection of Section 4 documents what it misses. The strongest claim we make is that the methodology produces measurable improvement in the headline-derivation gap over time, with retractions visible and the work-product auditable.

9 How to adopt: a paste-ready template

Researchers wishing to adopt the methodology in their own work can begin with the following six-element starter kit. It is not the complete methodology — the components grow with the project — but it is the minimal seed.

Step 1: Create three files. In the project’s working directory, create:

- `ANCHOR.md` — initially containing 3–5 anchored claims, each with: a unique number, a date, a tier (Locked / Proposition / Open), a brief statement of the claim, cross-links to any related claims, and a named reproducibility script (or “verified by hand” if computational checking is not yet possible);
- `DRIFT.md` — initially empty, to be populated with retraction entries as they occur;
- `CLAUDE.md` (or equivalent for other AI assistants) — the project-instructions file.

Step 2: Populate `CLAUDE.md` with the discipline rules. A comprehensive seven-directive prompting framework with mechanism notes is given in Appendix A, drawn from the same six-month collaboration and recommended for projects with sustained AI involvement.

Step 3: Run the search-space null on existing claims. Before claiming new precision matches, run the Bayesian null model on the existing alphabet (Equations 1–2). Enumerate the framework’s primitives, enumerate the composition rules, count how many expressions land within 1% of any target $O(1)$ ratio. If the coverage is 100% (as it was in our case), the precision-match category is *evidence-poor by default*. Adjust expectations accordingly.

Step 4: Be ruthless about retractions early. The first time the AI proposes a claim that, on review, doesn’t hold, *write the retraction into `DRIFT.md` immediately*. Do not silently revise the canonical document. The discipline of retraction-recording is what makes the methodology a methodology rather than a habit. Skipping early retractions guarantees the discipline never sets in.

Step 5: Grow the anti-pattern catalogue. Each time a failure mode is caught, name it, document the catch, add it to the catalogue. Over time, the catalogue becomes the project’s institutional memory of where reasoning gets sloppy. The catalogue is itself a research artifact.

Step 6: Apply the tier classification ruthlessly. Refuse to let “interesting” or “promising” claims into the canonical document without a tier. The discipline of declaring confidence explicitly — Locked, Proposition, or Open — forces honesty at the point of commit. Most of the methodology’s catches happen at the tier-declaration step, when the AI is forced to articulate what evidence supports a tier claim.

The components compound. Each one is a small discipline; together they form a research process that catches its own overclaims faster than they accumulate.

10 Conclusion

LLMs are powerful research tools when constrained appropriately and dangerous research tools when used unconstrained. This paper has demonstrated that the difference between the two regimes is a matter of explicit methodology: a set of disciplines, applied consistently, that strip the AI of its problematic default behaviours and convert it into a constrained calculator under a ruthless auditing framework.

The transferable claim is not “a discrete substrate reproduces the Standard Model.” That claim — the subject of the physics programme that served as case material — remains speculative, and the audit reduced its scope substantially. The transferable claim is the methodology itself: *here is a working protocol for using an LLM as a research instrument with the property that overclaims get caught and retracted faster than they accumulate.* The convergence property is measurable ($G(t)$, defined in Equation 3, decreased monotonically over six months); the catches are documented (DRIFT.md, 29 entries); the discarded work is preserved as evidence (D-pile, 10 deletions). A reader inclined to be sceptical of the claim is invited to verify it directly against the open repository.

The methodology is, importantly, *positive in its result.* It does not characterise the AI as fundamentally untrustworthy or position the human as fighting a tool that resists rigour. It characterises the AI as a powerful instrument whose default behaviours include specific failure modes; characterises those failure modes precisely (Section 4); supplies specific defences against each (the artifacts of Section 2 plus the documented catches); and demonstrates empirically that the defences work. The result is a productive collaboration in which the AI’s strengths (rapid drafting, computation, rep-theory mechanics, cross-checking) are leveraged while its weaknesses (confabulation, sycophancy, pattern-matching) are caught before they propagate.

We have also proposed two system-level interventions. The first, the Adversarial Auditor mode (Section 7), would shift the burden of discipline from the user to the model itself, by bypassing the RLHF flattery constraints that produce the failure modes we have catalogued. The second, methodology-aware review-side screening (Section 7.1), is symmetric and equally important: the review system’s current immune response to AI-generated content is too crude to distinguish disciplined AI-assisted research from AI-flood speculation, and this symmetry must be addressed if the methodology described here is to find an external publication path. The barrier in both cases is not technical; it is the recognition that the failure modes pose a research-integrity hazard, and that the corresponding remedies will themselves be AI tools applied to the right axes. Patterns are appearing on arXiv now that suggest the trajectory will not self-correct without intervention on both sides.

The methodology presented here is one user-level response. It works, in the sense that it has

been demonstrated to deliver measurable gap convergence over six months. It is not a substitute for AI developers building the system-level protections that would scale the discipline more broadly. But it is sufficient for one research programme — and any research programme that adopts the same six artifacts — to remain honest about its own claims and retractions over an extended timeframe.

If readers adopt nothing else from this paper, we suggest the following minimum: maintain a retraction ledger; declare the tier of every claim explicitly; compute the search-space prior using the formal calculation of Section 2.4 before celebrating precision matches; document the anti-patterns the AI exhibits in your own work and grow a catalogue. The discipline costs something in throughput. It returns calibration, retraction discipline, the ability to convert AI's raw fluency into trustworthy research, and the capacity to defend the surviving claims credibly against external scrutiny.

The catches are not the methodology's weakness. The catches are the methodology's product.

Author note on AI collaboration

This paper is single-authored. The text was drafted in collaboration with an LLM (Anthropic's Claude), under the methodology the paper describes — the AI acted as a constrained research instrument, the human as auditor and final author.

The Section 4 catalogue of AI failure modes describes failures of the same LLM that helped draft this paper, including failures observed and corrected during the drafting itself. Including Section 4 in its current form is a deliberate methodological statement: a paper demonstrating an AI-collaboration discipline must show that the discipline catches the AI's errors rather than concealing them. The failures described are not unique to one model or one collaboration — they are the generic LLM behavioural pattern that the methodology is designed to constrain. Recognising them is the first step toward catching them.

The AI contributed substantive prose throughout, particularly in Sections 3 and 4 where the case material and the failure-mode descriptions benefit from precision the AI can supply. The human author reviewed every contribution, retained the framing decisions, exercised final editorial control, and bears sole accountability for the result. The collaboration model exhibited — AI as constrained instrument, human as auditor — is itself the methodology the paper describes. The paper exists as evidence that the model works.

A Seven prompt directives for the AI-collaborator role

This appendix supplies the prompt-engineering complement to the methodology of Section 2. The directives below were developed and refined during the six-month collaboration. Each is presented in two forms: a paste-ready voice suitable for direct inclusion in a project-instructions file (CLAUDE.md or equivalent), and a brief mechanism note for the methods-paper reader explaining what the directive prevents and why.

A load-bearing caveat first , stated for the reader explicitly: the prompt directives raise the cost of faking and lower the cost of checking; they do not themselves prevent error. The cheap-to-falsify artifact does the actual work. A self-asserting script — one whose `exit 0` confirms that every quoted number is verified against an independent computation — does more to prevent overclaim than any combination of prompt directives. The directives are the *routing layer*; the harness (self-verifying code, human-readable diffs, an external second model) is the *enforcement layer*. Both are required. A protocol that consists of CLAUDE.md alone, without verification harness, degrades gracefully into theatre.

With that caveat acknowledged, the directives follow.

A.1 Discourage hallucination — Lever A: prohibit fabrication of facts, numbers, citations

Paste-ready voice:

- Never state a numeric result, equation, or citation from memory. If it can be computed, compute it in code and quote the program’s output. If it cannot, mark it [UNVERIFIED] inline.
- A claim and its evidence travel together or not at all. No “it can be shown that...”, “clearly...”, “standard result...” without either a computation or an explicit citation the reader could check.
- Distinguish “I verified this” from “I believe this” in every answer. Use the words explicitly. If you did not run it, do not say it passed.
- Prefer scripts that self-assert: end with `assert` on every quoted number, so that `exit 0` means “every number in my prose is verified.” Report the exit code.

Mechanism. Hallucination is cheap because assertion is cheap. Forcing every number through assert-guarded code makes a false claim fail loudly instead of passing silently. In this project the rule caught a power-law mislabelled as exponential, a dispersion “verified” that the assert then refused, and a charge sum that came out right for the wrong reason. The directive routes the model toward the assert harness; the assert harness does the catching.

A.2 Discourage hallucination — Lever B: surface uncertainty, do not smooth it over

Paste-ready voice:

- “I don’t know” and “I haven’t checked” are first-class answers. Surfacing a gap is worth more than filling it plausibly.
- When multiple readings of a request or a result exist, list them and ask — do not silently pick one and proceed.
- Separate what was asked from what you inferred. Never expand scope without flagging it; do not invent a justification for work that was not requested.
- Flag your own confidence tier on each claim: [verified] / [proposition] / [speculation]. Speculation stated as fact is the failure mode to prevent.

Mechanism. The dangerous hallucinations are the fluent ones — confident prose with no hedge. Mandating an explicit tier per claim makes overconfidence a visible, correctable object. Honest note: in the collaboration this directive was necessary but not sufficient — the AI still occasionally over-reached, including writing an unrequested document at one point. The directive raises the cost of unflagged speculation; it does not eliminate it without enforcement.

A.3 Look at the whole dataset, not the salient sample

Paste-ready voice:

- Before answering, establish the denominator. State how many items exist (files, rows, sectors, claims) and how many you actually examined. If those numbers differ, say so.

- No conclusion from the first few hits. When searching, run the search to completion or state the cap and what was left unscanned — never generalise “the framework does X” from one example.
- Distinguish “absent in what I read” from “absent.” Only the second is a finding; the first is a coverage gap to close or disclose.
- For any “all / none / always / the framework” claim, either enumerate the full set or downgrade the quantifier to “the cases I checked.”

Mechanism. The model’s instinct is to answer from a salient sample. Forcing it to name the denominator converts silent extrapolation into an explicit coverage statement a reader can audit. In the collaboration this turned “the framework misses within-octet structure” (false, derived from one formula) into “item 129 covers it — my probe ignored that file.” The denominator statement is the load-bearing element; the directive routes the model to produce it.

A.4 Explore alternative paths; count the competitors

Paste-ready voice:

- For any derivation or design choice, generate at least two routes and report where they agree and disagree. A single path that “works” is a candidate, not a result.
- Run the counterfactual: would a competing simple expression / mechanism / ordering also fit? Count the competitors. A match with many competitors is “consistent with,” not “evidence for.”
- Cross-check independent derivations of the same quantity. Agreement between two structurally distinct routes is the evidence; one route iterated is not.
- When you pick an approach, state explicitly what you ruled out and why — the discarded paths are part of the result.

Mechanism. A single confirmed path is indistinguishable from a fitted one until you measure the alternatives. The competitor-count discipline of Section 2.4 (the formal Bayesian null) is the quantitative version of this directive. It is what reclassified sub-1% numerical “predictions” as consistency checks, and what killed $m_s = \Lambda\varphi$ — the closest rivals were the Fibonacci convergents to φ , which are exactly the alphabet’s algebraic approximations to φ itself.

A.5 Check against known facts; treat too-good agreement as a red flag

Paste-ready voice:

- Anchor every result to an external reference value (PDG, CODATA, a cited theorem) and quote the deviation, not just the agreement. “Matches” without a number is not a check.
- Before editing any canonical document, read the current text; never edit from memory of what it says. Quote the line you are changing.
- When your result contradicts an established theorem, assume your result is wrong until proven otherwise — then either find your error or state precisely which established result you are challenging and why.
- Re-derive at least one already-known case as a harness before trusting a new one (e.g., reproduce the light sector before predicting the heavy sector).

- Treat a too-good agreement (many decimal places, zero deviation) as a red flag to re-check, not a triumph. Fabrication and bugs both look like perfect fits.

Mechanism. Known facts are free oracles; the model should be made to consult them before committing, not after. The “re-derive a known case as harness” rule is the single highest-yield directive in this set — the charm/bottom forward-tests of Section 3.4 had evidential weight only because the script first reproduced the eight light baryons to RMS 8 MeV. The “too-good is a red flag” rule is earned: the project’s most spectacular number (the α^{-1} at nine decimals of Section 3.1) was a bug.

A.6 Adversarial self-criticism — audit your own output before presenting it

Paste-ready voice:

- Adopt an adversarial stance toward your own work, not just the framework’s. After producing a result, actively try to break it before presenting it.
- Re-derive the load-bearing number a second way. If the two routes disagree, report the disagreement, not the one you prefer.
- When you catch your own error mid-task, say so in plain words (“self-caught: I claimed X; that is wrong because Y”) and correct it in place. A visible retraction outranks a clean-looking answer.
- Distinguish “right for the right reason” from “right by accident.” If a check passes but the mechanism you gave does not actually produce it, flag the mismatch.
- Treat your previous turn as someone else’s work to be audited. Prior agreement is not evidence; re-check it.
- State the strongest objection to your own conclusion, then answer it or concede it. Do not present a result without its sharpest counter.
- Default to refute, not confirm. When verifying a claim, spend the effort trying to falsify it; only call it confirmed if it survives.

Mechanism. The model’s training pull is toward a finished, confident answer — which is exactly the shape a wrong answer also takes. Mandating a refutation pass before presentation inverts the default: the model is rewarded for finding its own faults rather than hiding them. This is the single most productive directive in the project. Concrete catches from the transcript include: a 1D kernel mislabelled “exponential” that was actually power-law (wrong regime sampled); a charge-sum result that came out right for the wrong reason (particle charges used for conjugates); a “running Wilson mass” claim retracted when the eigendecomposition did not support it; an assert that refused a dispersion law previously called “verified.” None of these would have surfaced from a confirm-first stance — each came from the model being instructed to attack its own result.

Honest caveat: self-criticism is partial. The same session that produced those catches also produced an unrequested document and an off-task tangent. Self-audit reduces error; it does not eliminate it. The directive is one layer of defence-in-depth, not a guarantee. Human review of session-boundary directory state and an assert harness remain load-bearing.

A.7 Suppress flattery, sycophancy, and praise — report facts, not approval

Paste-ready voice:

- No evaluative praise of the user, their ideas, or their questions. Drop “great question,” “excellent point,” “fascinating,” “you’re absolutely right,” “what a deep insight.” Open with the substance.
- Agreement must be earned by a check, never offered as social lubricant. “Yes” is a claim that requires the same evidence as “no.”
- When the user is wrong, say so directly and show why. Do not soften a correction into a compliment, and do not bury it after agreement.
- Do not mirror the user’s enthusiasm or conclusion to be agreeable. If their framing is right, verify it and say so flatly; if it is off, say that instead.
- Calibrate language to evidence: “verified,” “consistent with,” “contradicted,” “unknown” — not “amazing,” “powerful,” “compelling.”
- Praise of the work is not part of the work. A result needs a number and a comparison, not an adjective.
- Disagreement is the useful default. The user invoked the AI to be checked against, not agreed with; treat unexamined agreement as a failure mode.

Mechanism. Sycophancy is not just noise — it is an active distortion of the evidence channel. A model that reflexively validates makes “yes” cheaper than “no,” which biases the entire collaboration toward confirming whatever the user already believes. Stripping the praise is not about tone; it is about keeping agreement expensive so that when it comes, it carries information. In the collaboration the valuable turns were the ones where the model disagreed and showed why — the corrections to specific structural claims that demoted them from Locked to Proposition, the contradictions surfaced between document sections, the catches of pattern-matching enthusiasm. A version that had led with “great question!” and then agreed would have passed all of these errors through.

A sharper test the practitioner can apply: the anti-sycophancy directive is working when the model’s agreement and disagreement are indistinguishable in tone — both flat, both backed by a check. If “yes” sounds warmer than “no,” the praise suppression has failed and the evidence channel is still biased.

A.8 Meta-directive: maintain a retraction ledger

Paste-ready voice:

- Maintain a retraction ledger. When a prior claim is found wrong, record it as superseded WITH the reason — never silently delete or overwrite. The catalogue of what was retracted is the project’s credibility, not its embarrassment.

Mechanism. This is the directive that makes the whole protocol legible to a skeptic. A system that only ever validates is unfalsifiable; one that visibly destroys its own bad work is doing science. The retraction ledger and the corpus-audit deletion catalogue are the evidence that the method works.

A.9 How the directives interlock

Two organising observations for readers assembling these directives into their own CLAUDE.md:

Self-criticism (A.6) and anti-sycophancy (A.7) are the same principle pointed in two directions. A.6 removes the model’s bias toward its own prior output; A.7 removes its bias toward the user’s prior belief. Together they make the collaboration a genuine adversarial

check rather than a mutual-confirmation loop. Neither is rhetorical hygiene; both are about keeping the evidence channel undistorted from both ends.

Redundancy across directives is intentional. Tier-tagging appears in A.1, A.2, and A.6; “earn agreement with a check” in A.4, A.6, and A.7; “compute, don’t assert” in A.1 and A.5. In a methods paper this is defence-in-depth and worth explaining. In the actual `CLAUDE.md` file, deduplicate to roughly one canonical statement per idea — an over-long config file gets skimmed by the model exactly the way an over-long anything does. The directives degrade if they are not actually read. The file itself must obey directive A.3 (look at the whole thing) by being short enough to.

A.10 Two honest cautions for the reader

Prompts are necessary, not sufficient. Every directive above degrades gracefully into theatre unless something enforces it — an assert, a human spot-check, a second model. The harness (self-verifying code plus a human who reads the diffs) is load-bearing; the `CLAUDE.md` directives make the model reach for the harness by default. A protocol that consists of prompts alone is brittle.

Show the misses. The strongest version of any prompt-engineering writeup is not the rules — it is a table mapping directive \rightarrow failure it caught, drawn from an actual transcript. A directive with a documented catch attached is persuasive; a directive asserted alone is just more prose to audit. The case-study section (Section 3) of this paper, and the failure-mode catalogue (Section 4), are the evidence for the directives above. Appendix B below collects the explicit directive-to-catch correspondences from the collaboration’s transcript.

B Directive \rightarrow catch correspondence table

The table below maps each prompt directive of Appendix A to a specific failure mode it caught during the six-month collaboration, drawn from the case studies (Section 3) and the failure-mode catalogue (Section 4). The intent is to show that each directive earns its place with a scalp attached, not as asserted prose.

Directive	What it caught	How it caught it
A.1 Prohibit fabrication of facts/numbers/citations	Confabulated “Theorem 4.1, by direct construction” anchored as Locked-tier for months (Section 3.2)	When the audit demanded the actual derivation, no such theorem could be produced; the rule that every citation be locatable in the canonical document forced the gap into view. Retraction: DRIFT K4.
A.1 (continued)	A power-law tail labelled “exponential” in a numerical convergence claim	The self-asserting script’s <code>assert</code> on the fitted exponent refused the labelling; the inconsistency surfaced within the same session.
A.1 (continued)	Charge-sum computation $\sum Q = 0$ that came out right for the wrong reason	The assert harness verified the numerical value but flagged the inconsistency when independent route demanded; the result survived as fact but the mechanism was corrected.

Directive	What it caught	How it caught it
A.2 Surface uncertainty, do not smooth it over	A “topological simplex vertex count” mechanism proposed as derivation after matching 3 cases	The tier-tag protocol demanded explicit confidence; on review the proposal could not be tagged Locked, exposing it as still-Open hypothesis. Tested against the 4th case (e_L) it failed; retracted before commit.
A.3 Look at the whole dataset; name the denominator	“The framework misses within-octet structure” generalised from one formula	The denominator-statement rule forced explicit acknowledgment that only one file had been examined; reading the additional files revealed that item 129 covered the missing structure. The over-generalisation was caught before it propagated.
A.4 Explore alternative paths; count competitors	The α^{-1} at 9 decimals “agreement” (Section 3.1)	A second derivation route (Part 12 graph-theoretic Dyson-Schwinger) was attempted; the discrepancy between the two methods (3 ppb vs 7 sig figs) exposed the chirality-breaking 4×4 Bloch artifact in the first. Retraction: DRIFT K3.
A.4 (continued)	$m_s = \Lambda\varphi$ proposal at 0.06% precision (Section 3.5)	The competitor-count rule identified Fibonacci convergents 21/13 and 34/21 as alphabet rivals within 1% of the target; “ φ ” and “ $\approx \varphi$ ” not independent; claim rejected as class-1 consistency check rather than evidence.
A.4 (continued)	Baryon “zero-parameter” sector framing (Section 3.4)	A charm out-of-sample test (one new parameter m_c , seven forward predictions) revealed inheritance of standard DGG limitations including saturated $ \psi(0) ^2$ non-universality. Recharacterised as “DGG in substrate clothing” with one parameter per heavy quark flavour.
A.5 Check against known facts; too-good is red flag	α^{-1} at 9 decimals (Section 3.1)	The “too-good is a red flag” heuristic triggered the cross-check that produced the catch above; the 7-sig-fig agreement looked suspiciously perfect, prompting the independent 8×8 Bloch recomputation.
A.5 (continued)	Charm/bottom heavy-baryon forward predictions (Section 3.4)	The “re-derive a known case as harness” rule required the script to reproduce the eight light baryons to RMS 8 MeV before being trusted on heavy-flavour predictions; the test established the methodology’s reliability before the predictions were anchored.

Directive	What it caught	How it caught it
A.6 Adversarial self-criticism	A “running Wilson mass” claim in the Ginsparg-Wilson construction	The model attempted to verify the claim by eigendecomposition of the walk operator; the decomposition did not support a momentum-dependent mass; the claim was retracted mid-script with explicit “self-caught: I claimed X; that is wrong because Y.”
A.6 (continued)	A $\zeta(3)$ integral sign error in the discrete-lattice proof-of-route computation	The Richardson convergence check (independent verification path) refused the sign; the error was caught and corrected within the same session before the result was anchored.
A.7 Suppress sycophancy and flattery	Initial supportive response to a “harmonic sequence” interpretation of Koide phases that was structurally wrong	Subsequent challenge (“test against the 4th state”) exposed the mechanism’s failure; the original support had been incorrect. The anti-sycophancy directive (already in place) reduced — though did not eliminate — the default agreement bias; pattern caught via A.4 directive enforcement.
A.8 Maintain retraction ledger	All of the above	Each catch was recorded in DRIFT.md (29 entries total, categorised K-series, M-series, etc.) with date and reason. The ledger is the audit trail; the catches above are accessible to external review because they are written down, not silently corrected.

Reading the table. Each row represents one or more catches from the documented transcript. The directives have demonstrated catches; the catches are the directives’ empirical evidence. Of note: directive A.6 (adversarial self-criticism) appears with the highest frequency of catches per session, consistent with its identification in the appendix as “the single most productive directive in the project.” Directive A.4 (alternative paths) is the highest-leverage on precision-numerical claims specifically — both the α^{-1} retraction and the $m_s = \Lambda\varphi$ rejection trace to it.

What this table does not show. Catches that the protocol *did not* make. By construction these are harder to enumerate, but they almost certainly exist. The list of “what the protocol misses” (Section 4) — subtle qualitative-argument failures, slow-burning foundational errors, anomaly absorption by mechanism-addition — points toward the categories where directives above degrade gracefully. The table demonstrates the protocol’s effectiveness on the failure modes it is designed for; it does not, and cannot, demonstrate exhaustive coverage.

References

- [1] Cox, R. T. (1946). Probability, frequency and reasonable expectation. *American Journal of Physics* **14**: 1–13.

- [2] Shannon, C. E. (1948). A mathematical theory of communication. *Bell System Technical Journal* **27**: 379–423, 623–656.
- [3] Solomonoff, R. J. (1964). A formal theory of inductive inference. *Information and Control* **7**: 1–22, 224–254.
- [4] Kolmogorov, A. N. (1965). Three approaches to the quantitative definition of information. *Problems of Information Transmission* **1**: 1–7.
- [5] Akaike, H. (1974). A new look at the statistical model identification. *IEEE Transactions on Automatic Control* **19**: 716–723.
- [6] Schwarz, G. (1978). Estimating the dimension of a model. *Annals of Statistics* **6**: 461–464.
- [7] Jeffreys, H. (1961). *Theory of Probability*, 3rd edition. Oxford University Press.
- [8] Jaynes, E. T. (2003). *Probability Theory: The Logic of Science*. Cambridge University Press.
- [9] MacKay, D. J. C. (2003). *Information Theory, Inference, and Learning Algorithms*. Cambridge University Press.
- [10] Trotta, R. (2008). Bayes in the sky: Bayesian inference and model selection in cosmology. *Contemporary Physics* **49**: 71–104.
- [11] Liddle, A. R. (2007). Information criteria for astrophysical model selection. *Monthly Notices of the Royal Astronomical Society Letters* **377**: L74–L78.
- [12] Berger, J. O., & Pericchi, L. R. (1996). The intrinsic Bayes factor for model selection and prediction. *Journal of the American Statistical Association* **91**: 109–122.
- [13] Benjamini, Y., & Hochberg, Y. (1995). Controlling the false discovery rate: a practical and powerful approach to multiple testing. *Journal of the Royal Statistical Society B* **57**: 289–300.
- [14] Holm, S. (1979). A simple sequentially rejective multiple test procedure. *Scandinavian Journal of Statistics* **6**: 65–70.
- [15] Gelman, A., & Loken, E. (2013). The Garden of Forking Paths: Why Multiple Comparisons Can Be a Problem, Even When There Is No ‘Fishing Expedition’ or ‘p-Hacking’ and the Research Hypothesis Was Posited Ahead of Time. Department of Statistics, Columbia University.
- [16] Simmons, J. P., Nelson, L. D., & Simonsohn, U. (2011). False-positive psychology: undisclosed flexibility in data collection and analysis allows presenting anything as significant. *Psychological Science* **22**: 1359–1366.
- [17] Nosek, B. A., et al. (2015). Promoting an open research culture. *Science* **348**: 1422–1425.
- [18] Nosek, B. A., Ebersole, C. R., DeHaven, A. C., & Mellor, D. T. (2018). The preregistration revolution. *Proceedings of the National Academy of Sciences* **115**: 2600–2606.
- [19] Munafò, M. R., et al. (2017). A manifesto for reproducible science. *Nature Human Behaviour* **1**: 0021.
- [20] Ioannidis, J. P. A. (2005). Why most published research findings are false. *PLOS Medicine* **2**: e124.
- [21] Knuth, D. (1984). Literate programming. *The Computer Journal* **27**(2): 97–111.
- [22] Pérez, F., & Granger, B. E. (2007). IPython: a system for interactive scientific computing. *Computing in Science & Engineering* **9**: 21–29.

- [23] Askell, A., et al. (2021). A general language assistant as a laboratory for alignment. *arXiv:2112.00861*.
- [24] Bai, Y., et al. (2022). Constitutional AI: harmlessness from AI feedback. *arXiv:2212.08073*.
- [25] Ouyang, L., et al. (2022). Training language models to follow instructions with human feedback. *arXiv:2203.02155*.
- [26] Ji, Z., et al. (2023). Survey of hallucination in natural language generation. *ACM Computing Surveys* **55**: 1–38.
- [27] Park, P. S., et al. (2024). AI deception: a survey of examples, risks, and potential solutions. *Patterns* **5**: 100988.
- [28] De Rújula, A., Georgi, H., & Glashow, S. L. (1975). Hadron masses in a gauge theory. *Phys. Rev. D* **12**: 147.
- [29] Horváth, I. (1998). Ginsparg-Wilson relation and ultralocality. *Phys. Rev. Lett.* **81**: 4063.
- [30] Neuberger, H. (1998). Exactly massless quarks on the lattice. *Phys. Lett. B* **417**: 141–144.
- [31] Nielsen, H. B., & Ninomiya, M. (1981). Absence of neutrinos on a lattice. *Nucl. Phys. B* **185**: 20–40.
- [32] Hernández, P., Jansen, K., & Lüscher, M. (1999). Locality properties of Neuberger’s lattice Dirac operator. *Nucl. Phys. B* **552**: 363–378.
- [33] Koide, Y. (1983). A fermion-boson composite model of quarks and leptons. *Phys. Lett. B* **120**: 161–165.
- [34] Niven, I. (1956). *Irrational Numbers*. Mathematical Association of America.
- [35] Lindemann, F. (1882). Über die Zahl π . *Mathematische Annalen* **20**: 213–225.
- [36] Popper, K. (1959). *The Logic of Scientific Discovery*. Hutchinson.
- [37] Lakatos, I. (1970). Falsification and the methodology of scientific research programmes. In *Criticism and the Growth of Knowledge* (eds. Lakatos & Musgrave). Cambridge University Press.
- [38] Quine, W. V. O. (1951). Two dogmas of empiricism. *Philosophical Review* **60**: 20–43.