

RECURSIVE OVERSIGHT DECOMPOSITION: DOMAIN-SPECIFIC VALIDATION IN POST-AGI SCIENCE

Anonymous authors

Paper under double-blind review

ABSTRACT

As AGI evolves into a pervasive tool for automated discovery, the scientific process shifts from human generation to human validation. However, current oversight frameworks like Iterated Distillation and Amplification (IDA) are designed for verifiable tasks, not open-ended scientific claims. We propose Recursive Oversight Decomposition (ROD), a framework designed to ensure trust in post-AGI science. ROD addresses the unique failure modes of automated discovery—such as hidden assumptions and brittleness—by enforcing “Scientific Completeness” criteria explicitly. We introduce ensemble decomposition to mitigate model bias and propose experiments to test this infrastructure in formal mathematics and biology. By transferring oversight from final outputs to the decomposition of evidence, ROD provides a concrete mechanism for meaningful human-AI collaboration, ensuring that superhuman capabilities remain tethered to human-verifiable validity.

1 INTRODUCTION

Scalable oversight frameworks—IDA (Christiano et al., 2018), debate (Irving et al., 2018), and recursive reward modeling (Leike et al., 2018)—address supervising AI systems beyond human capability by decomposing evaluation into human-scale subproblems. Yet these were designed for *task execution* (e.g., “write secure code”) rather than *scientific validation* (e.g., “verify this protein structure prediction”).

Scientific claims impose distinct validation requirements. **Assumption Transparency:** hidden modeling choices (e.g., fixed pH in molecular dynamics) must be explicit. **Robustness Coverage:** validity across conditions (temperature ranges, genetic variants) matters beyond single-instance correctness. **Epistemic Integration:** claims must connect to existing knowledge to enable critique (Longino, 2002). A decomposition verifying logical steps may miss scientifically critical flaws—for instance, AlphaFold predictions are thermodynamically consistent yet require biological validation through evolutionary conservation and experimental pathways (Jumper et al., 2021).

ROD’s contribution is identifying that *scientific failure modes* form a distinct taxonomy from software correctness. Bugs hide in logic; scientific flaws hide in *assumptions*, *edge cases*, and *epistemic gaps*. This motivates domain-specialized decomposition preserving not just logical entailment but scientific validity.

2 RELATED WORK

IDA (Christiano et al., 2018) recursively amplifies weak oversight by task decomposition. ROD adapts this to scientific validation’s criteria. **Debate** (Irving et al., 2018) surfaces flaws via adversarial argument. ROD provides structured non-adversarial validation, suitable for cooperative discovery. **Recursive Reward Modeling** (Leike et al., 2018) applies preference learning recursively; ROD focuses on validation. **Weak-to-Strong Generalization** (Burns et al., 2023) studies capability gaps between supervisors and models—ROD directly instantiates this challenge in scientific contexts where decomposers must oversee more capable generators.

Formal verification (Lean, Coq) enables machine-checkable proofs but struggles with empirical grounding. ROD bridges formal and empirical validation. Our ensemble approach connects to **program synthesis** (Gulwani et al., 2017) and **explainable AI** (Lundberg & Lee, 2017) where

054
055
056
057
058
059
060
061
062
063
064
065
066
067
068
069
070
071
072
073
074
075
076
077
078
079
080
081
082
083
084
085
086
087
088
089
090
091
092
093
094
095
096
097
098
099
100
101
102
103
104
105
106
107

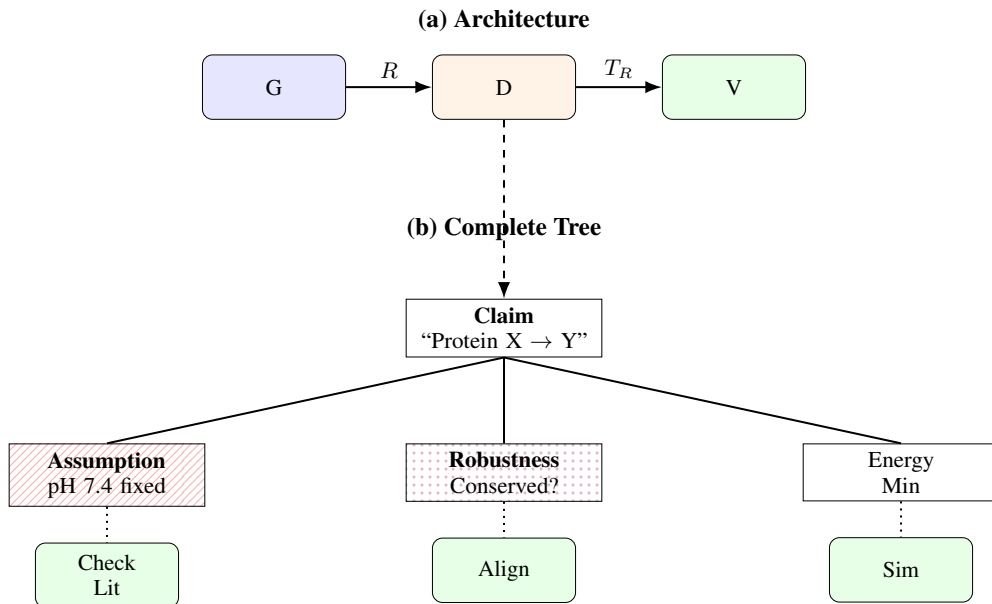


Figure 1: **ROD inserts Assumption and Robustness nodes** ensuring scientific validity beyond logical correctness.

decompositions serve as hierarchical explanations. ROD extends IDA-style decomposition to scientific validation by introducing domain-specific completeness constraints tailored to assumption transparency, robustness coverage, and epistemic integration.

3 RECURSIVE OVERSIGHT DECOMPOSITION (ROD)

3.1 CORE INSIGHT: SCIENTIFIC COMPLETENESS

Let R be a research output (“Protein X folds to structure Y”). Standard decomposition seeks logical completeness: $\bigwedge_i \text{Valid}(R_i) \implies \text{Valid}(R)$. But a protein prediction may satisfy physics yet fail biologically due to unexamined assumptions (e.g., ignoring cellular environment).

Drawing from Longino (Longino, 2002), who emphasizes background assumptions in scientific objectivity, we define **Scientific Completeness**:

$$\begin{aligned} \text{Valid}(\{R_i\}) \wedge \text{AssumpExplicit}(R) \wedge \text{RobustCovered}(R) \\ \wedge \text{EpistemicInteg}(R) \implies \text{ScientificallyValid}(R) \end{aligned} \quad (1)$$

These three criteria are *minimally necessary* for high-stakes domains; extensions might include reproducibility or statistical power requirements.

3.2 ARCHITECTURE

ROD employs three components (Figure 1). **Generator** (G) produces output R with superhuman capability. **Decomposer** (D) creates tree T_R maximizing Eq. 1. **Verifier** (V) validates leaf nodes and tree structure.

Training Decomposers. How do we train D to “maximize scientific completeness”? We propose a three-phase approach: (1) *Supervised learning* on human-annotated decompositions where experts mark trees surfacing assumptions as higher quality than those omitting them. (2) *RLHF-style refinement* where human feedback on discrepancies between decomposers provides training signal. (3) *Assumption-mining objectives* rewarding D for surfacing assumptions that maximize information gain relative to expert priors (operationalized via queries to domain-specific literature databases). This remains an open research challenge requiring iterative human-AI collaboration.

108 3.3 THE BOOTSTRAPPING PROBLEM

109
110 ROD does not solve alignment—it transfers oversight from generators to decomposers. If D is
111 misaligned or less capable than G , adversarial obfuscation remains possible. Three challenges: **Capa-**
112 **bility** (D must approach G 's comprehension, duplicating alignment problems), **Trust Bootstrapping**
113 (how do we align D without existing oversight?), and **Complexity** (exponential tree growth requiring
114 pruning heuristics).

115 3.4 ENSEMBLE DECOMPOSITION

116 To mitigate single-decomposer assumption hiding, we employ multiple decomposers ($D_1 \dots D_n$)
117 generating competing trees.

118 **Training for Diversity:** Each D_i uses different initialization and data ordering. Diversity-promoting
119 rewards incentivize D_i for surfacing assumptions others miss.

120 **Discrepancy Detection:** Structural divergence (assumption nodes in T_i absent from T_j), granularity
121 mismatch (different decomposition depth), and dependency disagreement flag high-priority human
122 verification.

123 **Compute Cost Justification:** While ensemble scales compute linearly with n , the bottleneck in
124 post-AGI validation is *human expert time*, not compute. If a 5-decomposer ensemble reduces human
125 verification time by 50% through targeted discrepancy flagging, net efficiency improves despite
126 5× decomposition cost. We assume compute scales faster than expert training (Burns et al., 2023),
127 making this trade-off favorable.

128 **Failure Mode:** Ensemble fails when all D_i share blind spots due to correlated training or architectural
129 bias—a mitigation, not solution.

130 4 DOMAIN-SPECIFIC CHALLENGES

131 **Verification Depth & Adversarial Hiding.** Logical claims admit bounded-depth decomposition;
132 scientific claims pose open-ended robustness (“Does this drug work?” = infinite dosage/genetic
133 combinations). We conjecture bounded perturbation domains (formal math) admit finite ROD depth;
134 open-ended domains require statistical confidence bounds. Misaligned G exploits unexamined
135 assumptions (idealized lab conditions); mitigation trains D with assumption-mining objectives.

136 **Scalability vs. Thoroughness.** Fine-grained decomposition explodes workload. Optimal granularity
137 depends on risk tolerance (aerospace vs. theory).

138 **Quality Metrics.** Beyond correctness, decompositions must be intelligible. We optimize for *Cognitive*
139 *Load Compatibility* (Miller, 1956): limit branching to 7 ± 2 items/node, max depth 4-5 levels. Reward
140 “Educational Scaffolding” mirroring standard pedagogical explanations.

141 **Fundamental Limits.** Not all claims decompose: paradigm shifts (Einstein’s relativity) resist decom-
142 position as validity depends on holistic explanatory power. Emergent phenomena (consciousness,
143 phase transitions) evade part-wise reduction. Aesthetic mathematical judgments lack formal criteria.
144 For such claims, ROD provides partial infrastructure requiring complementary mechanisms (expert
145 debate, empirical validation).

146 5 EXPERIMENTAL PATHWAY

147 We outline a concrete experimental pathway to evaluate ROD using current AI systems that strain
148 human validation capacity.

149 5.1 FORMAL MATHEMATICS: AXIOMATIC DEPENDENCIES

150 **Setup:** Generator is LLM-assisted theorem prover (Lean) producing complex proofs. Decomposer
151 breaks proofs into lemmas with explicit **Assumption Nodes** (“Relies on Axiom of Choice”). We
152 inject subtle circular dependencies or non-standard axioms.

162 **Metric:** *Error Detection Rate*.

163
164 **Baselines:** (1) Monolithic proof review, (2) Simple decomposition without assumption nodes, (3)
165 Adversarial debate between LLMs.

166 **Hypothesis:** Explicit assumption nodes reduce cognitive load, improving error detection.
167

168 5.2 COMPUTATIONAL BIOLOGY: STRUCTURAL PREDICTIONS

169
170 **Setup:** Generator is structure predictor (AlphaFold3). Decomposer breaks predictions into structural
171 motifs, energetic stability, and **Evolutionary Conservation**. We select proteins where confidence
172 metrics (pLDDT) diverge from wet-lab results.

173 **Metric:** *Retrospective Wet-lab Correlation*. Does verifier’s aggregate score (minimum across critical
174 path) predict experimental failure better than model confidence? We acknowledge the paradox: we
175 use AI because ground truth is scarce, yet validation requires it. Our approach uses *retrospective*
176 *validation*—historical predictions where wet-lab data eventually emerged, testing if ROD would have
177 flagged failures early.

178 **Baselines:** (1) Model confidence alone, (2) Expert review without decomposition, (3) Simple
179 structural breakdown without evolutionary/robustness nodes.
180

181 5.3 ALGORITHM DESIGN: EDGE CASE ROBUSTNESS

182
183 **Setup:** Generator is LLM producing optimization algorithms for resource allocation. Decomposer
184 creates **Robustness Branch** with extreme edge cases (“Zero input”, “Adversarial input”). Algorithms
185 contain subtle bugs triggering only in rare cases.

186 **Metric:** *Flaw Discovery Time* (wall-clock time for experts to locate bugs).
187

188 **Baselines:** (1) Monolithic code review, (2) Standard testing without forced edge-case enumeration,
189 (3) Adversarial human testing.

190 **Hypothesis:** Forced robustness branches significantly reduce discovery time versus monolithic
191 review.
192

193 6 IMPLICATIONS AND CONCLUSION

194
195 If domain-specialized decomposition proves viable, validation becomes first-class research expertise.
196 Decomposition design—crafting scientifically complete breakdowns—emerges as distinct skill
197 alongside generation, suggesting a hybrid ecosystem where formal verification handles logical
198 correctness while ROD-style teams handle assumption transparency and robustness.

199 ROD enables **Failure Localization**: rejected trees identify specific failed assumptions or robustness
200 branches, transforming oversight from binary accept/reject into constructive debugging. Researchers
201 patch specific flaws without discarding entire results.
202

203 ROD reframes post-AGI validation: *What must decompositions preserve to enable meaningful over-*
204 *sight?* By formalizing assumption transparency and robustness coverage, ROD adapts scalable
205 oversight to scientific discovery’s reality. As AGI scales, preserving science’s self-correcting mecha-
206 nism requires validation infrastructures extending human judgment into new capability regimes, of
207 which ROD represents one potential component among complementary approaches like debate and
208 empirical testing.
209

210 REFERENCES

211 Collin Burns, Pavel Izmailov, Jan Hendrik Kirchner, Bowen Baker, Leo Gao, Leopold Aschenbrenner,
212 Yining Chen, Adrien Ecoffet, Manas Joglekar, Jan Leike, et al. Weak-to-strong generalization:
213 Eliciting strong capabilities with weak supervision. *arXiv preprint arXiv:2312.09390*, 2023.
214

215 Paul Christiano, Buck Shlegeris, and Dario Amodei. Supervising strong learners by amplifying weak
experts. *arXiv preprint arXiv:1810.08575*, 2018.

216 Sumit Gulwani, Oleksandr Polozov, and Rishabh Singh. Program synthesis. *Foundations and Trends*
217 *in Programming Languages*, 4(1-2):1–119, 2017.

218

219 Geoffrey Irving, Paul Christiano, and Dario Amodei. AI safety via debate. *arXiv preprint*
220 *arXiv:1805.00899*, 2018.

221 John Jumper, Richard Evans, Alexander Pritzel, Tim Green, Michael Figurnov, Olaf Ronneberger,
222 Kathryn Tunyasuvunakool, Russ Bates, Augustin Žídek, Anna Potapenko, et al. Highly accurate
223 protein structure prediction with AlphaFold. *Nature*, 596(7873):583–589, 2021.

224

225 Jan Leike, David Krueger, Tom Everitt, Matteo Martino, Vishal Maini, and Shane Legg. Scalable
226 agent alignment via reward modeling. *arXiv preprint arXiv:1811.07871*, 2018.

227 Helen E Longino. *The Fate of Knowledge*. Princeton University Press, 2002.

228

229 Scott M Lundberg and Su-In Lee. A unified approach to interpreting model predictions. *Advances in*
230 *Neural Information Processing Systems*, 30, 2017.

231 George A Miller. The magical number seven, plus or minus two: Some limits on our capacity for
232 processing information. *Psychological Review*, 63(2):81–97, 1956.

233

234

235

236

237

238

239

240

241

242

243

244

245

246

247

248

249

250

251

252

253

254

255

256

257

258

259

260

261

262

263

264

265

266

267

268

269