

---

# EDA Bench: An Execution-Based Benchmark for Language Model Agents on Functional PCB Construction

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 Language-model agents are increasingly able to operate electronic design automa-  
2 tion tools: they can search datasheets, select components, edit KiCad projects, run  
3 design checks, and iterate after failures. Existing evaluations, however, rarely test  
4 the property that matters most for printed-circuit-board reconstruction: whether  
5 the produced board preserves the external electrical behavior of the target de-  
6 sign. We introduce EDA Bench, an execution-based benchmark for agentic KiCad  
7 PCB reconstruction. EDA Bench contains 35 released tasks derived from public  
8 open-hardware projects spanning connector breakouts, power-path circuits, micro-  
9 controller systems, motor drivers, camera and compute adapters, and high-density  
10 carriers. Each harness receives a prompt and a standard KiCad environment, then  
11 must produce a complete KiCad project. Submissions are graded by an oracle that  
12 parses the realized PCB, extracts routed copper, checks declared external ports,  
13 models trace parasitics and short faults, runs ngspice simulations, and applies  
14 KiCad ERC and DRC checks. The score is based on observable external-I/O  
15 behavior rather than text overlap, screenshot similarity, component inventory, or  
16 reference-designator matching. We validate the artifact with reference, fail, and  
17 mutation canaries: all 35 released frozen references score 1.0, and all 35 released  
18 structural-failure canaries score at most 0.15. Initial baselines show that EDA  
19 Bench is not saturated by current evaluated harnesses. The strongest evaluated  
20 web-enabled harness scores 4.58% overall, building 21 of 35 projects while failing  
21 most functional checks. These results identify a substantial gap between producing  
22 syntactically valid or visually plausible KiCad artifacts and reconstructing boards  
23 that behave correctly at their electrical interfaces.

## 24 1 Introduction

25 Language-model agents can now operate development tools over long horizons. In software, this has  
26 made execution-based evaluation a central benchmark design pattern: an agent modifies an artifact,  
27 runs tools, observes failures, and is evaluated by executable checks rather than by surface-form  
28 similarity. SWE-bench evaluates agents on real GitHub issues whose solutions must pass repository  
29 tests [Jimenez et al., 2023]. OSWorld extends execution-based evaluation to open-ended desktop and  
30 web tasks in real computer environments [Xie et al., 2024]. ProgramBench evaluates a still broader  
31 reconstruction setting in which agents rebuild software from an executable and documentation, with  
32 behavioral tests used to compare against the reference program [Yang et al., 2026]. These benchmarks  
33 share a common principle: the evaluated artifact should be judged by behavior in an execution  
34 environment.

35 Hardware design has an analogous workflow. Engineers search for parts, inspect datasheets, draft  
36 schematics, route boards, run ERC and DRC checks, simulate where possible, and produce manu-  
37 facturing files. For electronic design automation, however, execution-based testing is substantially  
38 harder. Correctness depends on routed connectivity, electrical behavior, manufacturability constraints,  
39 and external interfaces rather than only on executable program behavior. Existing CAD and EDA  
40 benchmarks cover important parts of this space, including text-to-CAD generation [Khan et al., 2024],  
41 geometric CAD validation [Epoch AI, 2026, CAD Arena, 2026, Barkley et al., 2026], schematic  
42 generation [Al Hasan et al., 2026, Zou et al., 2026, Luo et al., 2026], schematic understanding  
43 [Lu et al., 2026], board-level schematic design [Qiu et al., 2026], hardware bug repair [Cui et al.,  
44 2026], and PCB placement and routing reasoning [Li et al., 2026]. These evaluations motivate  
45 execution-based assessment for engineering artifacts, but they do not directly test whether an agent  
46 can reconstruct a routed KiCad PCB that preserves the target board’s external electrical behavior.

47 EDA Bench addresses this gap by evaluating complete agent harnesses on source-backed KiCad  
48 PCB reconstruction tasks. A harness receives a fixed task prompt and runtime assets, then must  
49 write a complete KiCad project in `/workspace/final_project`. KiCad provides an open-source  
50 EDA environment for schematic capture, integrated simulation, PCB layout, 3D rendering, and  
51 manufacturing export [KiCad Project, 2026]. The released tasks are based on public hardware  
52 projects, and the default evaluated harnesses are web-enabled. This creates a tradeoff between  
53 reproducibility and realism similar to public software-engineering benchmarks: public sources make  
54 the benchmark auditable, but web-enabled agents may benefit from locating and adapting upstream  
55 material.

56 The benchmark is built around a functional scoring principle: evaluate the observable behavior of the  
57 realized board at its external electrical interfaces. A submitted project is parsed, checked, simulated,  
58 and compared against task contracts. ngspice supplies the open-source circuit-simulation backend  
59 used by the oracle [ngspice contributors, 2026]. A board can score well with a different internal  
60 implementation if it exposes the required external behavior. It can also fail with a visually plausible  
61 layout if external nets are missing, shorted, swapped, isolated, or not driven by the required circuitry.  
62 Although this evaluation requires more engineering effort than syntax checks, visual matching, or  
63 component-inventory comparison, positive scores are more informative because the submitted board  
64 must satisfy observable electrical-interface checks.

65 **Contributions.** This paper makes five contributions.

- 66 • We introduce EDA Bench, a public benchmark of 35 source-backed KiCad PCB reconstruc-  
67 tion tasks spanning connector breakouts, power-path circuits, controller boards, compute  
68 adapters, motor drivers, and high-density carrier boards.
- 69 • We define a reproducible harness protocol in which agents receive only task-visible runtime  
70 packs and must produce complete KiCad projects, while references, source snapshots,  
71 contracts, canaries, and oracle internals remain part of the grading pack.
- 72 • We develop an execution-based PCB scoring oracle that evaluates realized board behavior  
73 through external-I/O contracts, routed-copper extraction, ngspice simulation, task-family  
74 checks, active-circuit realization checks, KiCad ERC and DRC diagnostics, and manufac-  
75 turability caps.
- 76 • We validate the scoring artifact with frozen references, structural fail canaries, mutation  
77 canaries, and invariance controls, showing that the oracle accepts human-designed references  
78 while rejecting common function-breaking PCB failures.
- 79 • We report public-provenance baselines for six web-enabled harness configurations, including  
80 final projects, grader outputs, usage summaries, cost estimates, and regrade commands,  
81 and show that current evaluated frontier-model harnesses remain far from functional PCB  
82 reconstruction.

## 83 2 Benchmark overview

84 EDA Bench evaluates PCB reconstruction as an execution task. For each task, a harness receives  
85 a fixed prompt, runtime task pack, tool environment, model configuration, and access policy, then  
86 must produce a complete KiCad project in `/workspace/final_project`. The submitted project

Table 1: EDA Bench evaluation pipeline. Each submission is evaluated as a complete harness output rather than as text, a screenshot, or a component list.

Stage	Artifact or operation
Task input	The harness receives a fixed prompt, runtime task pack, tool environment, model configuration, and access policy.
Harness execution	The agent operates in the Dockerized KiCad environment and writes a project to <code>/workspace/final_project</code> .
Project parsing	The grader locates KiCad files and extracts connectors, pads, nets, board outline, and routed copper geometry.
Contract binding	External ports from the submitted board are matched to task-declared I/O contracts and expected measurements.
Execution checks	The grader applies KiCad ERC and DRC diagnostics, ngspice simulations, routed-I/O checks, parasitic models, and short-fault checks.
Score reporting	The oracle returns a normalized task score, hard-cap diagnostics, build status, and provenance artifacts.

87 is graded automatically. The oracle parses the KiCad files, extracts external connectors and routed  
 88 copper, applies task-declared stimuli and probes, runs simulation where applicable, checks ERC and  
 89 DRC diagnostics, and assigns a normalized score based on the realized board behavior.

90 Table 1 summarizes the execution and grading pipeline. The benchmark is source-backed. Each  
 91 released task is derived from a public hardware project and has a frozen reference design, but the  
 92 evaluated runtime pack does not mount the reference project, source snapshot, canaries, grader  
 93 contracts, or oracle implementation. This split makes the artifact auditable while preserving a clean  
 94 separation between agent-visible inputs and grader-only materials. Because the released tasks are  
 95 public-source and the default baselines are web-enabled, EDA Bench measures PCB reconstruction  
 96 under a declared public-source access policy rather than closed-book circuit invention.

97 Scores should be interpreted as properties of complete evaluated harnesses, not isolated model  
 98 weights. A harness may include prompting, web search, source discovery, KiCad tool use, repair  
 99 loops, and project-generation logic. This systems-level framing is intentional because practical EDA  
 100 automation depends on whether an agent can operate the full workflow, not only on whether a model  
 101 can emit plausible schematic or layout text in one shot.

102 The central scoring claim is functional rather than visual. A board can score well with a different  
 103 internal implementation if it exposes the required external electrical behavior under the task contract.  
 104 Conversely, a board can fail despite looking plausible if external nets are missing, isolated, shorted,  
 105 swapped, unrouted, or not driven by the required active circuitry. The score is therefore not a  
 106 fabrication certificate, regulatory review, thermal analysis, EMC validation, or vendor-accurate silicon  
 107 simulation. It is an execution-based measure of declared external-I/O behavior under the released  
 108 oracle.

### 109 3 Benchmark tasks

110 EDA Bench contains 35 released source-backed tasks. Each task is anchored to a public upstream  
 111 hardware project and a frozen reference design. The public task set covers six difficulty levels, shown  
 112 in Table 2. Five additional candidate tasks were quarantined from the release because their upstream  
 113 redistribution terms require permission or acquisition-script replacement.

114 The task count reflects the cost of producing score-supported hardware evaluation instances rather  
 115 than a preference for small benchmarks. Each released task requires a source design, a frozen KiCad  
 116 reference, runtime and grading pack separation, an external-I/O contract, canary submissions, license  
 117 triage, and oracle support. We therefore prioritize tasks with executable functional checks over a  
 118 larger collection of weakly specified or weakly checked prompts.

119 The task catalog includes compact breakouts and converters such as `usb_c_female_breakout`,  
 120 `mcp23017_breakout`, and `rs485_transceiver_breakout`. It also includes power and control  
 121 boards such as `bq24295_power_path_board`, `picopd`, and `robotont_driver_board`,  
 122 compute and camera or display adapters such as `cm4_csi_adapter`, `m2_pcie_adapter`,

Table 2: Task distribution and representative categories in EDA Bench.

Difficulty	Tasks	Representative task families
Very easy	2	Compact connector breakouts and small adapters
Easy	5	I/O breakouts, logic-level conversion, USB/UART basics
Medium	4	Power-path boards, serial transceivers, small controller boards
Hard	12	Compute adapters, sensor/camera boards, multi-connector carriers
Very hard	7	Dense mixed-signal boards and larger interface carriers
Extreme	5	Large baseboards and high-density compute boards

123 and `ov9281_dual_camera_board`, and large carrier boards such as `cm4_baseboard`,  
 124 `jetson_orin_baseboard`, and `zynq_som`.

125 Each task has two artifact forms. The runtime task pack is mounted for the evaluated agent and  
 126 contains the prompt, task metadata, and runtime assets. It does not include the reference project,  
 127 source snapshot, canaries, grader contracts, or oracle implementation. The full grading pack is  
 128 public and contains the reference project, upstream snapshot, canary submissions, license notes,  
 129 `io_contract.json`, `functional_contract.json`, and local tests. This split supports repro-  
 130 ducibility while avoiding direct mounting of the answer into the evaluated container.

## 131 4 Harness protocol

132 A harness is evaluated as a complete system. It receives the task prompt and must produce a KiCad  
 133 project in `/workspace/final_project`. Built-in harnesses run inside Docker with KiCad, ngspice,  
 134 Python, Node, internet access, and the configured model interface. The current evaluated harnesses  
 135 are Codex CLI with GPT-5.5 web search at low, medium, high, and xhigh reasoning, Pi with Gemini  
 136 3.1 Pro Preview web high, and Pi with DeepSeek V4 Pro web high.

137 This protocol evaluates system architecture, not only a model weight. A one-shot generator, a  
 138 tool-using agent, and a multi-agent harness can all target the same task interface. This distinction  
 139 matters because EDA work is iterative: the model may need to create files, run KiCad tools, inspect  
 140 failures, search for upstream sources, and repair invalid output.

141 Each run uploads provenance to Hugging Face. Reports include prompts, command lines, usage  
 142 summaries, grader outputs, final projects, and generated artifacts. Large Docker image uploads and  
 143 raw model transcripts can be disabled to minimize storage usage, but reports still retain byte-count  
 144 summaries for omitted transcript content.

## 145 5 Scoring oracle

146 The scorer evaluates a submitted KiCad project through explicit I/O simulation. Unsupported,  
 147 missing, or malformed submissions receive zero score. Supported tasks declare their external ports,  
 148 stimuli, probes, simulator models, expected measurements, and tolerances. For each submission, the  
 149 oracle performs seven steps. It locates and parses the submitted KiCad project, runs KiCad ERC  
 150 and DRC diagnostics, extracts connectors, pads, nets, and routed copper geometry, binds submitted  
 151 external ports to the task contract, constructs task-declared simulation probes and stimuli, computes  
 152 task-family measurements and hard caps, and returns a normalized score with diagnostic artifacts.

153 The primary score is a normalized task score in  $[0, 1]$ . A benchmark aggregate is a per-task difficulty-  
 154 weighted mean:

$$S = \frac{\sum_{i=1}^N w(d_i) s_i}{\sum_{i=1}^N w(d_i)}, \quad w(d_i) \in \{1.0, 1.5, 2.0, 3.0, 4.0, 5.0\}. \quad (1)$$

155 We also report the unweighted mean over tasks. Build success means that the submission produced a  
 156 KiCad project that the grader could parse and evaluate. It is not equivalent to functional success.

157 The main score comes from three observable properties of the submitted PCB: external connector  
 158 semantics, routed-copper behavior, and task-required active and manufacturing plausibility. The  
 159 oracle parses the project, extracts PCB geometry, applies deterministic ngspice stimuli at declared

Table 3: Representative task-family checks layered on the generic external-I/O oracle. These are calibrated bindings on measured oracle signals, not vendor-accurate silicon or channel models.

Task family	Declared behavior	Common hard caps
Connector breakouts and simple adapters	Connector pins and power/ground pads, continuity, shorts, routed coverage, pin/net semantic agreement	Swapped pins, split nets, isolated pads, power-ground shorts
Serial and level-shifting interfaces	USB/UART/RS-485/logic-side boundary ports, signal transfer families, active-device realization	Missing active path, weak external mapping, missing routed I/O
Power-path and charger boards	Input, battery/load, configuration, and status ports, power path presence, external safety	Missing active power role, shorts, missing component-function profile
Camera/display/compute carriers	High-density connector families and differential pairs, routed families, pair plausibility	Missing high-speed pair families, weak routing, shorts, split nets
Large mixed-signal baseboards	Multi-rail power, GPIO, debug, serial, high-speed, and peripheral interfaces	Incomplete connector semantics, missing active realization, shorts

160 external ports, samples boundary waveforms, and applies hard caps for failures such as missing routed  
 161 copper, swapped external pins, shorts, missing active circuitry, unpowered active devices, invalid high-  
 162 speed pair families, and gross fabrication implausibility. Exact reference designators, component  
 163 inventory, footprint similarity, outline similarity, and route-shape similarity are diagnostics, not  
 164 primary score signals.

165 Within each task, declared measurements produce normalized sub-scores for external-port mapping,  
 166 routed continuity, simulation behavior, and task-family realization checks. These sub-scores are  
 167 combined according to the task contract to form a soft functional score. Hard caps then impose upper  
 168 bounds for structural failures that make the board externally incorrect or physically implausible,  
 169 such as shorts, missing routed I/O, swapped connector pins, or absent active circuitry. The final task  
 170 score is the resulting capped functional score. This design allows partial credit for partially correct  
 171 external behavior while preventing submissions with severe structural faults from receiving high  
 172 scores through unrelated checks.

## 173 6 Artifact validation

174 Because EDA Bench uses a task-specific oracle, the release includes controls that calibrate score  
 175 interpretation. The pass canary is the frozen reference project and must score 1.0 under the same  
 176 oracle. The fail canary removes active internal circuitry when available. Passive boards fall back to  
 177 external boundary removal. The fail canary must score below the release threshold, currently 0.75.

178 **Oracle validity.** We validate the oracle as a contract-level evaluator, not as a vendor-accurate  
 179 electrical simulator. Its intended validity criterion is whether score changes track engineering-relevant  
 180 changes to declared external-I/O behavior. Positive controls test that human-designed reference boards  
 181 are realizable and score 1.0. Negative controls test that structural changes a PCB engineer would  
 182 regard as function-breaking receive low scores. Invariance controls test that benign representation  
 183 changes, such as reference-designator renaming or plausible route-width edits, do not dominate the  
 184 score.

185 A full local sweep on May 6, 2026 verified all 35 released task canaries: every pass canary scored  
 186 1.0, and every fail canary scored at most 0.15. The generated-mutation sweep graded 162 mutation  
 187 canaries across the original 40-task candidate set. All tasks passed, with every reference scoring  
 188 1.0 and every structural fail canary scoring at most 0.15. Mutations covered missing board outlines,  
 189 unrouted external nets, isolated external pads, power or signal pads tied to ground, and missing  
 190 high-speed-like routes for high-speed carrier tasks.

Table 4: Oracle-validity evidence. The controls calibrate the released score as an external-I/O contract evaluator rather than as a complete manufacturing or vendor-silicon review.

Evidence	What it supports
35/35 released frozen human-designed reference PCBs score 1.0	The tasks are solvable, and the oracle accepts real upstream designs under their declared contracts.
35/35 released structural fail canaries score at most 0.15	The score is not reducible to file existence, KiCad syntax, or component-shell matching.
162 candidate-set mutation canaries pass validation	The oracle is sensitive to shorts, unrouted nets, isolated pads, missing outlines, and missing high-speed-like routing.
Refdes-renamed USB-C board scores 1.0	The score does not depend on exact connector reference designators.
Plausible USB-C route-width edit remains at least 0.85	The score does not require exact route-shape or width matching when behavior remains plausible.
Tiny-trace USB-C control scores at most 0.35	Manufacturability caps affect scores for physically implausible routing.

Table 5: Complete web-enabled harness baselines. The primary score is difficulty-weighted, and the unweighted score is the mean task score.

Harness	Weighted	Unweighted	Builds
Codex GPT-5.5 web low	0.0124	0.0143	31/35
Codex GPT-5.5 web medium	0.0279	0.0314	35/35
Codex GPT-5.5 web high	0.0293	0.0229	11/35
Codex GPT-5.5 web xhigh	0.0458	0.0500	21/35
Pi Gemini 3.1 Pro Preview web high	0.0113	0.0086	1/35
Pi DeepSeek V4 Pro web high	0.0000	0.0000	0/35

191 These controls do not show that the oracle replaces expert board review or captures vendor-specific  
 192 silicon behavior. They show that, for the benchmark’s declared external-I/O contracts, the score  
 193 accepts human-designed reference boards, rejects common function-breaking PCB mutations, and is  
 194 invariant to several superficial implementation details.

## 195 7 Experiments

196 We evaluate six web-enabled harness configurations on the 35 released tasks. All runs use the same  
 197 task prompts, task packs, Docker image, and grading protocol. Table 5 reports a May 6, 2026  
 198 no-model regrade of saved final projects under the current oracle revision.

199 No evaluated harness achieves a high absolute score. The main empirical finding is that build success  
 200 substantially overestimates functional PCB reconstruction. The medium reasoning-level Codex  
 201 harness builds all 35 released projects but scores only 2.79%. The xhigh harness builds fewer projects  
 202 than medium but scores higher. These results indicate that executable KiCad output is not a sufficient  
 203 metric for EDA agents and that reasoning and tool-use choices affect the distribution of functional  
 204 failures.

205 Difficulty breakdowns further indicate that the benchmark is not saturated by easy tasks. The xhigh  
 206 harness reaches 13.75% on medium tasks and 7.08% on hard tasks, but only 2.00% on extreme  
 207 tasks. The gap grows as boards require larger external contracts, active circuitry, dense routing, or  
 208 high-speed interfaces.

209 **Statistical interpretation.** The current result table reports complete released-task sweeps, not  
 210 repeated independent sampling campaigns. We therefore do not claim statistically significant pairwise  
 211 model rankings. The claims supported by the table are qualitative and artifact-level: the benchmark  
 212 is not saturated by current evaluated harnesses, build success is not functional success, and access  
 213 policy materially affects the evaluated harnesses under the same task and grader revision.

Table 6: Completed no-web calibration runs. These rows are reported separately from the default web leaderboard because access policy is an evaluated condition.

Harness	Weighted	Unweighted	Builds
Codex GPT-5.5 no-web medium	0.0000	0.0000	0/35
Codex GPT-5.5 no-web high	0.0000	0.0000	0/35
Codex GPT-5.5 no-web xhigh	0.0000	0.0000	0/35
Codex GPT-5.4 no-web xhigh	0.0000	0.0000	0/35
Codex GPT-5.4-mini no-web medium	0.0000	0.0000	0/35
Codex GPT-5.4-mini no-web high	0.0000	0.0000	0/35

## 214 7.1 Access-policy calibration

215 EDA Bench reports access policy as part of the evaluated system. Six no-web Codex calibration  
 216 sweeps use the same released tasks and grader revision as the web-enabled runs, but disable web  
 217 access. All six scored zero and produced no buildable projects, indicating that the default prompts  
 218 and built-in Codex harnesses rely on public-source discovery as well as KiCad operation. These rows  
 219 are calibration runs rather than closed-book capability measurements for all possible harnesses, so  
 220 we report them separately from the web-enabled leaderboard.

## 221 8 Failure analysis

222 The results support three benchmark-level conclusions. First, current evaluated frontier-model  
 223 harnesses are far from solving functional PCB reconstruction: the best evaluated configuration scores  
 224 only 4.58% overall on the released task set. Second, syntactic validity is not sufficient: a harness  
 225 can produce buildable KiCad projects for every released task and still fail most functional checks.  
 226 Third, harness architecture matters: Codex GPT-5.5 configurations with different reasoning settings  
 227 produce different build and score tradeoffs, so EDA Bench evaluates complete systems rather than  
 228 base models alone.

229 Saved grading artifacts expose the main failure modes. Among buildable Codex GPT-5.5 xhigh  
 230 submissions, the most common caps were split required external-net continuity, isolated external-I/O  
 231 pads, wrong external-I/O pin/net mapping, missing or weak component-function realization, missing  
 232 or implausible high-speed pair geometry, and weak active-device power integrity. The all-building  
 233 Codex medium run shows a similar pattern: wrong external mapping, missing component-function  
 234 realization, split required external nets, same-layer trace shorts, and isolated external-I/O pads.

235 These caps correspond to concrete functional failures rather than parser failures. Submissions often  
 236 instantiate connector footprints with plausible labels, but leave required I/O pads isolated, split  
 237 required nets across disconnected copper, swap external pin semantics, omit active circuitry, or route  
 238 high-speed-like interfaces implausibly. These errors explain why KiCad parseability and component  
 239 presence are weak proxies for functional PCB reconstruction.

## 240 9 Related work

241 **Execution-based software-agent benchmarks.** Execution-based software benchmarks evaluate  
 242 agents by running the artifacts they produce. SWE-bench frames software engineering as repository-  
 243 level issue resolution, where an agent edits a real codebase and must pass tests associated with the  
 244 target issue [Jimenez et al., 2023]. OSWorld evaluates multimodal agents in real desktop and web  
 245 environments with custom execution-based evaluation scripts [Xie et al., 2024]. ProgramBench  
 246 studies a more reconstructive setting in which agents rebuild software from an executable and  
 247 documentation, and behavioral tests are generated to compare the candidate program with the  
 248 reference executable [Yang et al., 2026]. EDA Bench adopts this behavioral-evaluation perspective  
 249 but applies it to PCB reconstruction, where correctness depends on realized geometry, external  
 250 connectivity, and electrical-interface behavior rather than software tests alone.

251 **CAD generation and geometric evaluation.** Text2CAD introduced a large-scale text-to-parametric-  
 252 CAD dataset and model for generating CAD command sequences from natural-language prompts

253 [Khan et al., 2024]. CadEval and CAD Arena provide evaluation settings for text-to-CAD models  
254 using geometric validity, rendering, and prompt-level comparisons [Epoch AI, 2026, CAD Arena,  
255 2026]. CADSmith further emphasizes programmatic geometric validation for generated CadQuery  
256 models [Barkley et al., 2026]. These works motivate execution-based evaluation for engineering  
257 artifacts, but their primary correctness signals are geometric rather than electrical. PCB reconstruction  
258 requires connectivity, routing, short avoidance, active-circuit realization, and external-I/O behavior in  
259 addition to shape or layout plausibility.

260 **Circuit, schematic, and PCB benchmarks.** Recent EDA benchmarks study several adjacent tasks.  
261 CircuitLM and PCBSchemaGen generate or synthesize circuit schematics from natural-language  
262 prompts and use structural or constraint-based checks to improve schematic validity [Al Hasan et al.,  
263 2026, Zou et al., 2026]. SchGen studies PCB schematic generation using semantic-grounded code  
264 representations [Luo et al., 2026]. OmniSch evaluates multimodal schematic understanding and  
265 schematic-to-graph reasoning over real-world diagrams [Lu et al., 2026]. HWE-Bench evaluates  
266 board-level schematic design with datasheet context and hardware bug repair over repository-scale  
267 hardware projects [Qiu et al., 2026, Cui et al., 2026]. PCB-Bench evaluates LLM reasoning over  
268 placement, routing, and PCB design comprehension using text, images, and complete PCB projects  
269 [Li et al., 2026]. EDA Bench differs by requiring an agent to produce a complete routed KiCad PCB  
270 project and by scoring the realized board through external-I/O simulation and PCB-geometry checks  
271 rather than schematic-only validity, visual reasoning, or placement/routing question answering.

## 272 10 Limitations

273 EDA Bench is not a hidden-answer or closed-book benchmark. The released tasks are public-source,  
274 and upstream designs may be discovered through web search. This determines the supported claim:  
275 scores measure EDA task solving under fixed prompts, tools, task-pack revisions, provenance, and  
276 access policy, not closed-book invention.

277 The task set is intentionally small compared with text benchmarks. It has 35 released tasks because  
278 each task requires a real source design, a reference project, task contracts, canaries, license review,  
279 and grader support. Breadth should grow through additional source-backed packs or a sealed private  
280 split, not by adding weakly checked or specified tasks.

281 The oracle is functional but not a substitute for manufacturing review. It uses deterministic generic  
282 behavioral models. It does not prove electrical safety, regulatory compliance, high-speed protocol  
283 compliance, thermal reliability, power sequencing correctness, or production readiness. Richer  
284 per-task SPICE, IBIS, S-parameter, and field-solver models would improve fidelity. The current  
285 calibrated task-family checks are deliberately conservative: they tighten measured oracle signals for  
286 four representative board families, but they are not vendor-accurate silicon or channel simulations.

287 Task artifacts inherit mixed upstream hardware licenses. Five candidate tasks with unlicensed,  
288 upstream-terms-only, noncommercial, or otherwise restricted source snapshots are quarantined  
289 from the public release until they receive explicit redistribution permission or acquisition-script  
290 replacements.

## 291 11 Conclusion

292 EDA Bench evaluates whether language-model agents can reconstruct PCB designs that function  
293 at their electrical interfaces. It is built from real public hardware projects, scored by explicit I/O  
294 simulation, and evaluated as a full harness rather than a single text response. Initial results show  
295 that current systems often produce syntactically valid KiCad artifacts, but rarely produce boards that  
296 satisfy functional boundary contracts. This discrepancy is the central evaluation target of EDA Bench.

## 297 References

298 Khandakar Shakib Al Hasan, Syed Rifat Raiyan, Hasin Mahtab Alvee, and Wahid Sadik. CircuitLM:  
299 A Multi-Agent LLM-Aided Design Framework for Generating Circuit Schematics from Natural  
300 Language Prompts. arXiv:2601.04505, 2026. URL: <https://arxiv.org/abs/2601.04505>.

301 Jesse Barkley, Rumi Loghmani, and Amir Barati Farimani. CADSmith: Multi-Agent CAD Generation  
302 with Programmatic Geometric Validation. arXiv:2603.26512, 2026. URL: [https://arxiv.org/  
303 abs/2603.26512](https://arxiv.org/abs/2603.26512).

304 CAD Arena. CAD Arena: Open Benchmark for AI-Generated Parametric CAD. 2026. URL:  
305 <https://cadarena.dev/>.

306 Fan Cui, Hongyuan Hou, Zizhang Luo, Chenyun Yin, and Yun Liang. HWE-Bench: Benchmarking  
307 LLM Agents on Real-World Hardware Bug Repair Tasks. arXiv:2604.14709, 2026. URL:  
308 <https://arxiv.org/abs/2604.14709>.

309 Epoch AI. CadEval. 2026. URL: <https://epoch.ai/benchmarks/cad-eval>.

310 Carlos E. Jimenez, John Yang, Alexander Wettig, Shunyu Yao, Kexin Pei, Ofir Press, and  
311 Karthik Narasimhan. SWE-bench: Can Language Models Resolve Real-World GitHub Issues?  
312 arXiv:2310.06770, 2023. URL: <https://arxiv.org/abs/2310.06770>.

313 Mohammad Sadil Khan, Sankalp Sinha, Talha Uddin Sheikh, Didier Stricker, Sk Aziz Ali, and  
314 Muhammad Zeshan Afzal. Text2CAD: Generating Sequential CAD Models from Beginner-to-  
315 Expert Level Text Prompts. In *Advances in Neural Information Processing Systems*, volume 37,  
316 2024. URL: <https://arxiv.org/abs/2409.17106>.

317 KiCad Project. KiCad Documentation. 2026. URL: <https://docs.kicad.org/>.

318 Jindong Li, Lianrong Chen, Bin Yang, Jiadong Zhu, Ying Wang, Yuzhe Ma, and Menglin Yang. PCB-  
319 Bench: Benchmarking LLMs for Printed Circuit Board Placement and Routing. In *International  
320 Conference on Learning Representations*, 2026. URL: [https://openreview.net/forum?id=  
321 Q5QLu7XTWx](https://openreview.net/forum?id=Q5QLu7XTWx).

322 Taiting Lu, Kaiyuan Lin, Yuxin Tian, Yubo Wang, Muchuan Wang, Sharique Khatri, Akshit Kartik,  
323 Yixi Wang, Amey Santosh Rane, Yida Wang, Yifan Yang, Yi-Chao Chen, Yincheng Jin, and  
324 Mahanth Gowda. OmniSch: A Multimodal PCB Schematic Benchmark For Structured Diagram  
325 Visual Reasoning. arXiv:2604.00270, 2026. URL: <https://arxiv.org/abs/2604.00270>.

326 Qinpei Luo, Ruichun Ma, Xinyu Zhang, and Lili Qiu. SchGen: PCB Schematic Generation  
327 with Semantic-Grounded Code Representations. OpenReview preprint, 2026. URL: [https:  
328 //openreview.net/forum?id=TyWs6rWWHb](https://openreview.net/forum?id=TyWs6rWWHb).

329 ngspice contributors. ngspice, the open-source SPICE circuit simulator. 2026. URL: [https:  
330 //ngspice.sourceforge.io/](https://ngspice.sourceforge.io/).

331 Weibo Qiu, Yinhao Xiao, and Runyu Pan. HWE-Bench: Can Language Models Perform Board-level  
332 Schematic Designs? arXiv:2603.18102, 2026. URL: <https://arxiv.org/abs/2603.18102>.

333 Tianbao Xie, Danyang Zhang, Jixuan Chen, Xiaochuan Li, Siheng Zhao, Ruisheng Cao, Toh Jing  
334 Hua, Zhoujun Cheng, Dongchan Shin, Fangyu Lei, Yitao Liu, Yiheng Xu, Shuyan Zhou, Silvio  
335 Savarese, Caiming Xiong, Victor Zhong, and Tao Yu. OSWorld: Benchmarking Multimodal  
336 Agents for Open-Ended Tasks in Real Computer Environments. arXiv:2404.07972, 2024. URL:  
337 <https://arxiv.org/abs/2404.07972>.

338 John Yang, Kilian Lieret, Jeffrey Ma, Parth Thakkar, Dmitrii Pedchenko, Sten Sootla, Emily McMilin,  
339 Pengcheng Yin, Rui Hou, Gabriel Synnaeve, Diyi Yang, and Ofir Press. ProgramBench: Can  
340 Language Models Rebuild Programs From Scratch? arXiv:2605.03546, 2026. URL: [https:  
341 //arxiv.org/abs/2605.03546](https://arxiv.org/abs/2605.03546).

342 Huanghaohe Zou, Peng Han, Emad Nazerian, and Alex Q. Huang. PCBSchemaGen: Constraint-  
343 Guided Schematic Design via LLM for Printed Circuit Boards. arXiv:2602.00510, 2026. URL:  
344 <https://arxiv.org/abs/2602.00510>.

345 **A Artifact availability and provenance**

346 The source repository contains the runner and harness code, task-pack loader and grader support, pub-  
347 lic documentation, package tests, and the task-pack manifest. Python wheels and source distributions  
348 exclude concrete task instances, reference projects, upstream snapshots, and canary submissions. The  
349 Hugging Face task dataset contains runtime task packs, full grader task packs, prompts, references,  
350 upstream snapshots, canary submissions, task metadata, Croissant metadata, and license/provenance  
351 notes.

352 The current public task dataset is hosted at [https://huggingface.co/datasets/eda-](https://huggingface.co/datasets/eda-bench-neurips-2026/eda-bench-tasks)  
353 [bench-neurips-2026/eda-bench-tasks](https://huggingface.co/datasets/eda-bench-neurips-2026/eda-bench-tasks). The paper results use task-pack revision  
354 [c64b9169d6d8c0d73614e14c7050cf036aeb1559](https://huggingface.co/datasets/eda-bench-neurips-2026/eda-bench-tasks). The public provenance dataset is hosted at  
355 <https://huggingface.co/datasets/eda-bench-neurips-2026/eda-bench-provenance>.  
356 The public code release is hosted at <https://github.com/CAD-bench/eda-bench>.

357 Run provenance includes task prompts, command lines, usage summaries, grader outputs, final  
358 projects, and generated artifacts. Storage-light campaigns can omit full Docker image uploads  
359 and raw model transcripts, but retain final reports, final projects, grading artifacts, and byte-count  
360 summaries for omitted transcript content.

361 **B Validation controls**

Table 7: Reviewer-facing validation controls for score interpretation.

Control	Observed result
Frozen reference project	35/35 released references score 1.0
Structural fail canary	35/35 released fail canaries score at most 0.15, below the 0.75 release threshold
Generated mutation canaries	162 generated mutation canaries across the original 40-task candidate set pass validation
Calibrated task-family checks	USB-C, RS-485, BQ24295, and M.2/PCIe references score 1.0 and their mutation canaries remain below threshold
Refdes-invariant equivalent board	USB-C reference with renamed connector refdes scores 1.0
Missing routed copper	USB-C and NanoUPDI route-removal controls score 0.0
Plausible route-width edit	USB-C route-width edit remains at least 0.85
Implausibly tiny traces	Manufacturability cap applies and USB-C tiny-trace control scores at most 0.35
Missing or unparsable project	Returns 0.0 and build failure

362 Two additional generic mutation families are implemented experimentally for future validation:  
363 swapped external connector pins and signal-bearing active devices with power/ground pads discon-  
364 nected. These should be staged in release packs or reported as validated only after a full canary  
365 sweep.

366 **C Failure-mode diagnostics**

367 **D Baseline provenance**

368 The May 6, 2026 result table regrades saved final projects without rerunning model har-  
369 nesses. Source runs were: `codex__gpt-5.5-web-low_20260505T200208Z_2719115`, `codex_`  
370 `_gpt-5.5-web-medium_20260506T010715Z_2812175`, `codex__gpt-5.5-web-medium_`  
371 `20260506T010715Z_2812182`, `codex__gpt-5.5-web-medium_20260506T032737Z_2884928`,  
372 `codex__gpt-5.5-web-medium_20260506T010716Z_2812190`, `codex__gpt-5.5-web-high_`  
373 `20260505T200208Z_2719113`, `codex__gpt-5.5-web-xhigh_20260506T032737Z_2884940`,  
374 `codex__gpt-5.5-web-xhigh_20260506T032737Z_2884938`, `codex__gpt-5.5-web-xhigh_`  
375 `20260506T032737Z_2884943`, `codex__gpt-5.5-web-xhigh_20260506T034126Z_2891808`,

Table 8: Common functional failure modes exposed by the oracle. These failures are typical of buildable submissions and illustrate why KiCad parseability is not sufficient for functional PCB reconstruction.

Failure mode	Observable symptom	Why it affects score
Split required external nets	Pins that should belong to the same external interface are assigned to disconnected copper regions or separate nets	Boundary stimuli cannot propagate through the declared interface, even when connector footprints are present
Isolated external I/O pads	Required pads are placed but not connected to routed copper or active circuitry	The board exposes a physical pad without realizing the electrical behavior required by the contract
Wrong pin or net mapping	External connector pins are swapped, mislabeled, or attached to semantically incorrect nets	The board may appear visually plausible while driving or observing the wrong external signal
Missing active realization	Components required for level shifting, power-path behavior, transceiver behavior, or control logic are absent or not connected	Passive connectivity alone cannot implement the task-required circuit behavior
Weak active-device power integrity	Active devices are present but lack required power, ground, or enable connectivity	The simulated or inferred behavior cannot be attributed to a realizable powered circuit
Implausible high-speed geometry	Differential or high-speed-like interface families are missing, weakly routed, or routed with implausible geometry	Dense carrier and compute tasks require not only pin labels but routed interface families with plausible physical realization

376 pi\_\_gemini-3.1-pro-preview-web-high\_20260505T184850Z\_2697190, and  
 377 pi\_\_deepseek\_\_deepseek-v4-pro-web-high\_20260505T200208Z\_2719116.

378 The executable regrade entry point is `uv run regrade-provenance`. For example, the Codex  
 379 medium row is reproduced by passing its four source runs:

```
380 uv run regrade-provenance \  

381   --run codex__gpt-5.5-web-medium_20260506T010715Z_2812175 \  

382   --run codex__gpt-5.5-web-medium_20260506T010715Z_2812182 \  

383   --run codex__gpt-5.5-web-medium_20260506T032737Z_2884928 \  

384   --run codex__gpt-5.5-web-medium_20260506T010716Z_2812190
```

385 By default the command uses the current released task manifest, downloads saved  
 386 submission/final\_project directories from Hugging Face provenance, and runs the current  
 387 local oracle.

388 Completed no-web calibration runs are present under these provenance paths:  
 389 evals/harnesses/codex\_\_gpt-5.5-no\_web-medium\_20260507T033319Z\_3432012,  
 390 evals/harnesses/codex\_\_gpt-5.5-no\_web-high\_20260507T033319Z\_3432018,  
 391 evals/harnesses/codex\_\_gpt-5.5-no\_web-xhigh\_20260507T033319Z\_3432017,  
 392 evals/harnesses/codex\_\_gpt-5.4-no\_web-xhigh\_20260507T033528Z\_3439753, evals/  
 393 harnesses/codex\_\_gpt-5.4-mini-no\_web-medium\_20260507T033724Z\_3448333, and  
 394 evals/harnesses/codex\_\_gpt-5.4-mini-no\_web-high\_20260507T033724Z\_3448341.

## 395 E Cost and compute accounting

396 The completed GPT-5.5 web runs used storage-light provenance. The combined Codex GPT-5.5  
 397 low/medium/high/xhigh campaign uploaded about 185.6 MiB of Hugging Face provenance. Across  
 398 all six final harnesses, the total uploaded provenance was about 193.1 MiB.

399 The estimates are computed from recorded token usage, current provider pricing, and the Codex  
 400 token-based rate card. They are included to make leaderboard submissions auditable and to discourage  
 401 hidden cost and provenance blowups.

Table 9: Recorded Codex GPT-5.5 web campaign accounting.

Harness	Estimated OpenAI API USD	Estimated Codex credits	HF provenance
Codex low	38.22	955.41	48.9 MB
Codex medium	84.40	2110.05	49.5 MB
Codex high	49.39	1234.64	27.6 MB
Codex xhigh	182.75	4568.64	68.6 MB

402 **F Release and license gate**

403 EDA Bench separates runtime prompts from full task packs with references, snapshots, canaries, and  
 404 grader contracts. For NeurIPS release, full packs with unlicensed, upstream-terms-only, noncom-  
 405 mercial, or otherwise restricted source snapshots must be handled by one of three policies: obtain  
 406 explicit redistribution permission, remove the source snapshot and provide a reproducible acquisition  
 407 script, or quarantine the task from the public full pack while retaining runtime metadata.

408 The current public release quarantines five candidate tasks under this conservative  
 409 rule: `cm5io_official`, `mixed_signal_stm32_dev_board`, `rp2350b_dev_board`,  
 410 `stm32f_audio_codec`, and `ef28_badge`. The public score-supported full-pack release has  
 411 35 tasks. These task ids may remain documented as candidates, but their references, source snapshots,  
 412 and canaries are not part of the released full-pack artifact.

413 **NeurIPS Paper Checklist**

414 **1. Claims**

415 Question: Do the main claims made in the abstract and introduction accurately reflect the  
 416 paper’s contributions and scope?

417 Answer: [Yes]

418 Justification: The abstract and introduction state the paper’s scope as introducing an  
 419 execution-based benchmark for web-enabled PCB reconstruction agents, documenting  
 420 its harness protocol and external-I/O oracle, validating the scoring artifact with canaries,  
 421 and reporting initial harness baselines. The paper explicitly excludes hidden-answer circuit  
 422 invention, fabrication readiness, regulatory certification, and vendor-accurate electrical  
 423 validation.

424 Guidelines:

- 425 • The answer [N/A] means that the abstract and introduction do not include the claims  
 426 made in the paper.
- 427 • The abstract and/or introduction should clearly state the claims made, including the  
 428 contributions made in the paper and important assumptions and limitations. A [No] or  
 429 [N/A] answer to this question will not be perceived well by the reviewers.
- 430 • The claims made should match theoretical and experimental results, and reflect how  
 431 much the results can be expected to generalize to other settings.
- 432 • It is fine to include aspirational goals as motivation as long as it is clear that these goals  
 433 are not attained by the paper.

434 **2. Limitations**

435 Question: Does the paper discuss the limitations of the work performed by the authors?

436 Answer: [Yes]

437 Justification: The Limitations section discusses public-source task exposure, task-set size,  
 438 oracle fidelity, and upstream license constraints. The benchmark overview and access-policy  
 439 calibration sections also clarify that the default baselines are web-enabled and should not be  
 440 interpreted as closed-book PCB invention.

441 Guidelines:

- 442 • The answer [N/A] means that the paper has no limitation while the answer [No] means  
 443 that the paper has limitations, but those are not discussed in the paper.

- 444 • The authors are encouraged to create a separate “Limitations” section in their paper.
- 445 • The paper should point out any strong assumptions and how robust the results are to
- 446 violations of these assumptions.
- 447 • The authors should reflect on the scope of the claims made, e.g., if the approach was
- 448 only tested on a few datasets or with a few runs.
- 449 • The authors should reflect on the factors that influence the performance of the approach.
- 450 • The authors should discuss the computational efficiency of the proposed algorithms
- 451 and how they scale with dataset size.
- 452 • If applicable, the authors should discuss possible limitations of their approach to
- 453 address problems of privacy and fairness.
- 454 • While the authors might fear that complete honesty about limitations might be used by
- 455 reviewers as grounds for rejection, a worse outcome might be that reviewers discover
- 456 limitations that are not acknowledged in the paper.

### 457 3. Theory assumptions and proofs

458 Question: For each theoretical result, does the paper provide the full set of assumptions and  
459 a complete (and correct) proof?

460 Answer: [N/A]

461 Justification: The paper is a benchmark and empirical evaluation paper and does not include  
462 theoretical results or proofs.

463 Guidelines:

- 464 • The answer [N/A] means that the paper does not include theoretical results.
- 465 • All the theorems, formulas, and proofs in the paper should be numbered and cross-
- 466 referenced.
- 467 • All assumptions should be clearly stated or referenced in the statement of any theorems.
- 468 • The proofs can either appear in the main paper or the supplemental material, but if
- 469 they appear in the supplemental material, the authors are encouraged to provide a short
- 470 proof sketch to provide intuition.
- 471 • Inversely, any informal proof provided in the core of the paper should be complemented
- 472 by formal proofs provided in appendix or supplemental material.
- 473 • Theorems and Lemmas that the proof relies upon should be properly referenced.

### 474 4. Experimental result reproducibility

475 Question: Does the paper fully disclose all the information needed to reproduce the main ex-  
476 perimental results of the paper to the extent that it affects the main claims and/or conclusions  
477 of the paper (regardless of whether the code and data are provided or not)?

478 Answer: [Yes]

479 Justification: The paper reports the task-pack revision, harness protocol, scoring formula,  
480 runtime and grading artifact split, provenance dataset, source run identifiers, and regrade  
481 procedure. The code, released task artifacts, grading artifacts, and saved final projects are  
482 public or described through the release and quarantine policy.

483 Guidelines:

- 484 • The answer [N/A] means that the paper does not include experiments.
- 485 • If the paper includes experiments, a [No] answer to this question will not be perceived
- 486 well by the reviewers: Making the paper reproducible is important, regardless of
- 487 whether the code and data are provided or not.
- 488 • If the contribution is a dataset and/or model, the authors should describe the steps taken
- 489 to make their results reproducible or verifiable.
- 490 • Depending on the contribution, reproducibility can be accomplished in various ways.
- 491 • While NeurIPS does not require releasing code, the conference does require all submis-
- 492 sions to provide some reasonable avenue for reproducibility, which may depend on the
- 493 nature of the contribution.

### 494 5. Open access to data and code

495 Question: Does the paper provide open access to the data and code, with sufficient instruc-  
496 tions to faithfully reproduce the main experimental results, as described in supplemental  
497 material?

498 Answer: [Yes]

499 Justification: Appendix A identifies the source repository, Hugging Face task dataset, task-  
500 pack revision, Croissant metadata, and public provenance dataset. Appendix F describes the  
501 license triage policy and states that five upstream-license-risk candidate tasks are quarantined  
502 from the released full-pack artifact.

503 Guidelines:

- 504 • The answer [N/A] means that paper does not include experiments requiring code.
- 505 • Please see the NeurIPS code and data submission guidelines ([https://neurips.cc/  
506 public/guides/CodeSubmissionPolicy](https://neurips.cc/public/guides/CodeSubmissionPolicy)) for more details.
- 507 • While we encourage the release of code and data, we understand that this might not be  
508 possible, so [No] is an acceptable answer.
- 509 • The instructions should contain the exact command and environment needed to run to  
510 reproduce the results.
- 511 • The authors should provide instructions on data access and preparation.
- 512 • The authors should provide scripts to reproduce all experimental results for the new  
513 proposed method and baselines.
- 514 • At submission time, to preserve anonymity, the authors should release anonymized  
515 versions, if applicable.
- 516 • Providing as much information as possible in supplemental material is recommended,  
517 but including URLs to data and code is permitted.

## 518 6. Experimental setting/details

519 Question: Does the paper specify all the training and test details (e.g., data splits, hyperpa-  
520 rameters, how they were chosen, type of optimizer) necessary to understand the results?

521 Answer: [Yes]

522 Justification: The benchmark overview, harness protocol, scoring oracle, artifact validation,  
523 and experiments sections describe the task interface, Dockerized tool environment, model  
524 and harness configurations, score definition, build-success definition, task count, validation  
525 controls, and access-policy separation.

526 Guidelines:

- 527 • The answer [N/A] means that the paper does not include experiments.
- 528 • The experimental setting should be presented in the core of the paper to a level of detail  
529 that is necessary to appreciate the results and make sense of them.
- 530 • The full details can be provided either with the code, in appendix, or as supplemental  
531 material.

## 532 7. Experiment statistical significance

533 Question: Does the paper report error bars suitably and correctly defined or other appropriate  
534 information about the statistical significance of the experiments?

535 Answer: [Yes]

536 Justification: The paper states that the reported results are complete released-task sweeps  
537 rather than repeated independent sampling campaigns. It therefore does not claim statistically  
538 significant pairwise model rankings. The supported conclusions are artifact-level and  
539 qualitative: the benchmark is not saturated by current evaluated harnesses, build success  
540 diverges from functional score, and access policy affects the evaluated harnesses under the  
541 reported protocol.

542 Guidelines:

- 543 • The answer [N/A] means that the paper does not include experiments.
- 544 • The authors should answer [Yes] if the results are accompanied by error bars, confidence  
545 intervals, or statistical significance tests, at least for the experiments that support the  
546 main claims of the paper.

- 547 • The factors of variability that the error bars are capturing should be clearly stated.
- 548 • The method for calculating the error bars should be explained.
- 549 • The assumptions made should be given.
- 550 • It should be clear whether the error bar is the standard deviation or the standard error
- 551 of the mean.
- 552 • For asymmetric distributions, the authors should be careful not to show symmetric
- 553 error bars that would yield results that are out of range.
- 554 • If error bars are reported in tables or plots, the authors should explain how they were
- 555 calculated and reference the corresponding figures or tables.

## 556 8. Experiments compute resources

557 Question: For each experiment, does the paper provide sufficient information on the com-  
558 puter resources (type of compute workers, memory, time of execution) needed to reproduce  
559 the experiments?

560 Answer: [Yes]

561 Justification: The paper specifies the Dockerized KiCad, ngspice, Python, and Node en-  
562 vironment, records run identifiers and provenance, and reports cost and provenance-size  
563 accounting for the main Codex campaign in Appendix E. Saved run provenance includes  
564 additional execution details for reproduced runs.

565 Guidelines:

- 566 • The answer [N/A] means that the paper does not include experiments.
- 567 • The paper should indicate the type of compute workers CPU or GPU, internal cluster,  
568 or cloud provider, including relevant memory and storage.
- 569 • The paper should provide the amount of compute required for each of the individual  
570 experimental runs as well as estimate the total compute.
- 571 • The paper should disclose whether the full research project required more compute  
572 than the experiments reported in the paper.

## 573 9. Code of ethics

574 Question: Does the research conducted in the paper conform, in every respect, with the  
575 NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines?>

576 Answer: [Yes]

577 Justification: The work is a benchmark and evaluation study without human subjects or  
578 deception. The paper discusses limitations, release constraints, and the risk of overtrust in  
579 generated hardware artifacts.

580 Guidelines:

- 581 • The answer [N/A] means that the authors have not reviewed the NeurIPS Code of  
582 Ethics.
- 583 • If the authors answer [No], they should explain the special circumstances that require a  
584 deviation from the Code of Ethics.
- 585 • The authors should make sure to preserve anonymity.

## 586 10. Broader impacts

587 Question: Does the paper discuss both potential positive societal impacts and negative  
588 societal impacts of the work performed?

589 Answer: [Yes]

590 Justification: The paper frames the benchmark as a tool for measuring EDA-agent reliability  
591 and explicitly states that scores are not evidence of manufacturing readiness, safety, EMC  
592 compliance, regulatory compliance, or production readiness. The main negative impact is  
593 overtrust in generated hardware artifacts. The mitigation is conservative scoring, validation  
594 controls, and explicit scope limits.

595 Guidelines:

- 596 • The answer [N/A] means that there is no societal impact of the work performed.

- 597
- 598
- 599
- 600
- 601
- 602
- 603
- 604
- 605
- 606
- 607
- If the authors answer [N/A] or [No], they should explain why their work has no societal impact or why the paper does not address societal impact.
  - Examples of negative societal impacts include potential malicious or unintended uses, fairness considerations, privacy considerations, and security considerations.
  - The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments.
  - The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, when it is being used as intended but gives incorrect results, and from misuse.
  - If there are negative societal impacts, the authors could also discuss possible mitigation strategies.

## 608 11. Safeguards

609 Question: Does the paper describe safeguards that have been put in place for responsible  
610 release of data or models that have a high risk for misuse (e.g., pre-trained language models,  
611 image generators, or scraped datasets)?

612 Answer: [N/A]

613 Justification: The paper releases benchmark tasks, grading assets, and evaluation code rather  
614 than a new high-risk generative model or scraped personal or sensitive dataset. Upstream-  
615 license-risk hardware snapshots are quarantined from the released full-pack artifact until  
616 permission or acquisition-script replacement is available.

617 Guidelines:

- 618
- 619
- 620
- 621
- 622
- 623
- 624
- The answer [N/A] means that the paper poses no such risks.
  - Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model.
  - Datasets that have been scraped from the Internet could pose safety risks.
  - We recognize that providing effective safeguards is challenging, and many papers do not require this, but we encourage authors to take this into account and make a best faith effort.

## 625 12. Licenses for existing assets

626 Question: Are the creators or original owners of assets (e.g., code, data, models), used in  
627 the paper, properly credited and are the license and terms of use explicitly mentioned and  
628 properly respected?

629 Answer: [Yes]

630 Justification: Appendix F states the license triage policy and identifies the five quarantined  
631 candidate tasks. The repository includes LICENSES.md, task-pack license notes, and source  
632 provenance metadata.

633 Guidelines:

- 634
- 635
- 636
- 637
- 638
- 639
- 640
- 641
- 642
- 643
- 644
- 645
- 646
- The answer [N/A] means that the paper does not use existing assets.
  - The authors should cite the original paper that produced the code package or dataset.
  - The authors should state which version of the asset is used and, if possible, include a URL.
  - The name of the license should be included for each asset.
  - For scraped data from a particular source, the copyright and terms of service of that source should be provided.
  - If assets are released, the license, copyright information, and terms of use in the package should be provided.
  - For existing datasets that are re-packaged, both the original license and the license of the derived asset should be provided.
  - If this information is not available online, the authors are encouraged to reach out to the asset's creators.

## 647 13. New assets

648 Question: Are new assets introduced in the paper well documented and is the documentation  
649 provided alongside the assets?

650 Answer: [Yes]

651 Justification: The paper documents the task taxonomy, runtime and full-pack split, scoring  
652 oracle, validation controls, provenance structure, Croissant metadata, artifact availability, and  
653 release policy. The repository and dataset contain detailed task and release documentation.

654 Guidelines:

- 655 • The answer [N/A] means that the paper does not release new assets.
- 656 • Researchers should communicate the details of the dataset/code/model as part of their  
657 submissions via structured templates.
- 658 • The paper should discuss whether and how consent was obtained from people whose  
659 asset is used.
- 660 • At submission time, remember to anonymize your assets if applicable.

#### 661 14. Crowdsourcing and research with human subjects

662 Question: For crowdsourcing experiments and research with human subjects, does the paper  
663 include the full text of instructions given to participants and screenshots, if applicable, as  
664 well as details about compensation (if any)?

665 Answer: [N/A]

666 Justification: The paper does not involve crowdsourcing or research with human subjects.

667 Guidelines:

- 668 • The answer [N/A] means that the paper does not involve crowdsourcing nor research  
669 with human subjects.
- 670 • Including this information in the supplemental material is fine, but if the main contribu-  
671 tion of the paper involves human subjects, then as much detail as possible should be  
672 included in the main paper.
- 673 • According to the NeurIPS Code of Ethics, workers involved in data collection, curation,  
674 or other labor should be paid at least the minimum wage in the country of the data  
675 collector.

#### 676 15. Institutional review board (IRB) approvals or equivalent for research with human 677 subjects

678 Question: Does the paper describe potential risks incurred by study participants, whether  
679 such risks were disclosed to the subjects, and whether Institutional Review Board (IRB)  
680 approvals (or an equivalent approval/review based on the requirements of your country or  
681 institution) were obtained?

682 Answer: [N/A]

683 Justification: The paper does not involve human-subjects research.

684 Guidelines:

- 685 • The answer [N/A] means that the paper does not involve crowdsourcing nor research  
686 with human subjects.
- 687 • Depending on the country in which research is conducted, IRB approval or equivalent  
688 may be required for any human subjects research.
- 689 • We recognize that the procedures for this may vary significantly between institutions  
690 and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the  
691 guidelines for their institution.
- 692 • For initial submissions, do not include any information that would break anonymity, if  
693 applicable, such as the institution conducting the review.

#### 694 16. Declaration of LLM usage

695 Question: Does the paper describe the usage of LLMs if it is an important, original, or  
696 non-standard component of the core methods in this research? Note that if the LLM is used  
697 only for writing, editing, or formatting purposes and does *not* impact the core methodology,  
698 scientific rigor, or originality of the research, declaration is not required.

699

Answer: [N/A]

700

Justification: LLMs are not important, original, or non-standard components of the core method development in this research.

701

702

Guidelines:

703

- The answer [N/A] means that the core method development in this research does not involve LLMs as any important, original, or non-standard components.

704

705

- Please refer to our LLM policy in the NeurIPS handbook for what should or should not be described.

706