
CAD-bench: Benchmarking Language Models on Functional CAD Generation

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Language-model agents are increasingly able to operate computer-aided design
2 (CAD) toolchains, but current evaluations for their capabilities in this domain
3 often measure only whether the models build, their rendered appearance, or crude
4 geometric similarity. These proxy and approximate criteria tend to be much easier
5 for the agents to meet, but fail to measure the properties that matter in mechan-
6 ical design–mating interfaces, exact dimensions, and functional behavior–while
7 overweighting less significant implementation details, such as wall thickness for an
8 outer covering around a gearbox. We introduce CAD-bench, an execution-based
9 benchmark for evaluating language-model CAD agents. CAD-bench contains 17
10 tasks across four difficulty tiers, ranging from basic solids to threaded mating
11 pairs and functional gear trains. Each submission is executed, exported as ge-
12 ometry, and evaluated with task-specific checks, including dimensional and pose
13 verification, thread-profile analyses, and Blender-based rigid-body simulation for
14 complex machines. The benchmark supports both one-shot CAD-code generation
15 and agent harnesses that produce final STEP artifacts in a Linux environment.
16 Initial results show that CAD-bench is not saturated by current systems. The best
17 standalone model reaches 59.9% overall, and functional tasks remain near zero for
18 most standalone runs. Agent harnesses perform better than one-shot generation,
19 but still fail frequently on the hardest tasks, reaching a maximum score of 68%. No
20 harness, whether one-shot or agent, can achieve a score of 50% on the insane-level
21 tasks. CAD-bench therefore demonstrates the current limited capability of models
22 to complete end-to-end CAD tasks, as well as the gap between measuring their
23 capabilities with functional evaluation and simpler methods.

24 1 Introduction

25 Computer-aided design (CAD) is a central representation for physical product development. CAD
26 artifacts encode not only visible shape, but also dimensions, constraints, features, interfaces, and
27 assemblies that downstream engineering tools use for simulation, inspection, and manufacturing.
28 As language-model agents become more capable of writing code, operating tools, and iterating on
29 execution feedback, CAD generation is becoming a natural potential target for agentic automation.

30 Evaluation, however, remains a bottleneck. A generated CAD artifact can compile, render, and
31 resemble a target object while still failing the engineering task. Common errors include shifted
32 poses, incorrect hole placement, missing mating features, invalid thread profiles, gear collisions,
33 intersecting bodies, and in general, assemblies that fail to transmit motion. These errors impact
34 whether a generated artifact is usable at all.

35 Prior CAD-generation benchmarks have established important foundations. SketchGraphs repre-
36 sents CAD sketches as geometric constraint graphs [Seff et al., 2020]. CAD-as-Language and

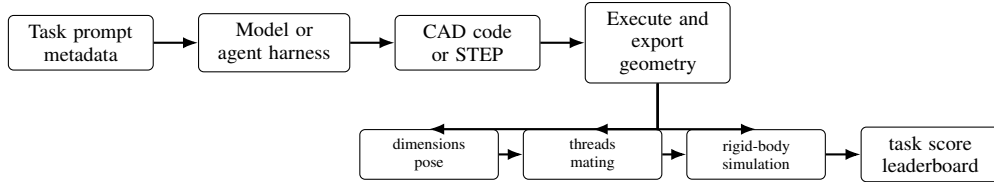


Figure 1: CAD-bench evaluation pipeline. A harness receives a task specification and produces CAD code or a STEP artifact. The runtime executes the submission, exports geometry, and applies task-specific checks for dimensions, pose, thread observables, mating behavior, and rigid-body function before aggregation.

37 DeepCAD model CAD sketches or construction histories as serialized design programs [Ganin
 38 et al., 2021, Wu et al., 2021]. Fusion 360 Gallery provides human-authored design sequences and a
 39 programmatic reconstruction environment [Willis et al., 2021]. More recent systems move toward
 40 natural-language CAD generation, executable CAD code, or geometric validation. Text2CAD studies
 41 text-to-parametric-CAD generation from natural-language prompts [Khan et al., 2024]. CADPrompt
 42 pairs natural-language CAD prompts with expert CAD code and uses a visual verification loop
 43 [Alrashedy et al., 2024]. CAD-Coder and Text-to-CadQuery target text-to-CAD-code generation
 44 with executable feedback or geometric rewards [Guan et al., 2025, Xie and Ju, 2025]. CadEval, CAD-
 45 Smith, and CAD Arena provide related evaluation settings based on rendering, geometry, validity,
 46 or programmatic measurement [Epoch AI, 2026, Barkley et al., 2026, CAD Arena, 2026]. These
 47 benchmarks and systems show that executable geometric feedback is valuable. CAD-bench, however,
 48 targets a complementary and more significant question: whether a submitted CAD artifact satisfies
 49 mechanically meaningful task requirements.

50 CAD-bench evaluates complete harnesses rather than a single modeling API or decoding format. A
 51 harness receives a task prompt and must produce either CAD code or a STEP artifact. The benchmark
 52 then executes the submission, exports geometry where needed, and applies task-specific evaluators.
 53 The current release includes 17 tasks across four tiers: basic solids, feature-rich parts, standards-like
 54 mechanical components, and functional assemblies. The hardest tasks require not only primitives and
 55 simple shapes but also interface compatibility and rigid-body function.

56 The central claim of CAD-bench is that build success is not a sufficient CAD metric. For example, in
 57 the reported agent rows, the right-angle gearbox builds in 96.9% of attempts but averages only 4.0%
 58 on the full functional score. CAD-bench is designed to distinguish artifacts that merely execute from
 59 artifacts that satisfy the dimensional, interface, and functional constraints of the task.

60 **Contributions.** This paper makes four contributions.

- 61 • We introduce CAD-bench, a 17-task executable benchmark for language-model CAD agents,
 62 spanning primitive solids, feature-rich parts, standards-like components, threaded mating
 63 pairs, and functional transmissions.
- 64 • We define a harness protocol that supports both one-shot CAD-code generation and agent-
 65 produced STEP submissions, allowing model-only and tool-using systems to be evaluated
 66 through a shared artifact interface.
- 67 • We develop a task-specific scoring stack that combines build success, dimensional and pose
 68 checks, reference-geometry gates, thread-profile analysis, and Blender-based rigid-body
 69 simulation.
- 70 • We report complete baseline sweeps for standalone models and agent harnesses, together
 71 with scorer-validation cases, aggregation-sensitivity checks, repeated-run diagnostics, and
 72 reproducibility artifacts.

73 2 Related Work

74 CAD-bench builds on CAD generation, engineering-document understanding, and executable agent
 75 benchmarks.

76 **CAD generation and CAD code benchmarks.** Seff et al. [2020] introduce SketchGraphs, a large
77 collection of CAD sketches represented as geometric constraint graphs. Ganin et al. [2021] formulate
78 CAD sketch generation as a language modeling problem over serialized design programs, while
79 Wu et al. [2021] model CAD construction histories as operation sequences. Willis et al. [2021]
80 provide human-authored Fusion 360 design sequences and an environment for programmatic CAD
81 reconstruction. More recently, Khan et al. [2024] introduce Text2CAD for generating sequential
82 CAD models from natural-language prompts, and Alrashedy et al. [2024] introduce CADPrompt,
83 a benchmark of natural-language CAD prompts paired with expert CAD code, along with a visual
84 verification loop. CAD-Coder [Guan et al., 2025] and Text-to-CadQuery [Xie and Ju, 2025] move
85 directly toward text-to-parametric-code generation, using executable CadQuery and geometric rewards
86 such as Chamfer distance. CadEval [Epoch AI, 2026] evaluates text-to-CAD systems with rendering
87 success, Chamfer and Hausdorff distances, and volume or consistency checks. CADSmith [Barkley
88 et al., 2026] is especially close in spirit: it uses execution repair plus programmatic OpenCASCADE
89 measurements and visual judging to refine CadQuery code. CAD Arena [CAD Arena, 2026] provides
90 an emerging public comparison grid for 20 prompts, but its current reported metric is valid executable
91 STL production.

92 These works establish CAD as a structured generation problem and show that executable geometric
93 feedback matters. CAD-bench differs by scoring submitted artifacts against task-level mechanical
94 requirements, including mating interfaces, thread observables, exact pose, and rigid-body function,
95 rather than construction-sequence, rendering, or point-cloud similarity.

96 **Engineering-document understanding.** Doris et al. [2024] introduce DesignQA for evaluating
97 multimodal understanding of engineering documentation. That work is complementary to CAD-
98 bench: understanding drawings and manuals is part of the input side of engineering intelligence,
99 while CAD-bench evaluates output-side artifact generation in a runnable CAD environment.

100 **Executable benchmarks for agents.** Recent work has shown the value of realistic, environment-
101 grounded evaluation for agents. AgentBench [Liu et al., 2023], GAIA [Mialon et al., 2023], Visu-
102 alWebArena [Koh et al., 2024], and OSWorld [Xie et al., 2024] benchmark agents in interactive
103 environments rather than static text tasks. In software engineering, SWE-bench [Jimenez et al., 2023]
104 and MLE-bench [Chan et al., 2024] show that end-to-end execution environments uncover failure
105 modes that synthetic benchmarks miss. RE-Bench [Wijk et al., 2024] and PaperBench [Starace et al.,
106 2025] extend this perspective to open-ended research engineering. CAD-bench adopts the same
107 benchmark philosophy for CAD: models should be evaluated in the artifact-producing environment
108 itself.

109 **Benchmark methodology.** BetterBench [Reuel et al., 2024] emphasizes best practices such as
110 transparent reporting, replicability, and careful benchmark lifecycle management, while LiveBench
111 [White et al., 2024] highlights the importance of contamination-resistant evaluation. CAD-bench
112 follows these principles through executable scoring, explicit artifact traces, and a design that makes
113 future task refreshes and extensions straightforward.

114 **3 Motivation: Why CAD Needs an Executable Benchmark**

115 CAD evaluation requires stronger checks than syntactic validity or visual plausibility. A model
116 can generate valid Python, export a watertight mesh, or render an object with the correct silhouette
117 while still producing the wrong dimensions, pose, feature placement, mating geometry, or assembly
118 behavior. In mechanical design, these local errors are often decisive.

119 Surface-level geometric similarity is also incomplete. Many CAD tasks are defined by constrained
120 interfaces rather than by global shape alone. A screw head with an incorrect socket, a thread that
121 cannot mate, a gear train that collides, or a gearbox that fails to transmit motion may remain visually
122 plausible under coarse rendering or point-cloud metrics. Such artifacts may look close to a reference
123 but fail the engineering role specified by the prompt.

124 CAD agents also operate in executable tool environments, where they must synthesize parametric
125 programs, reason over APIs and exchange formats, and satisfy exact toolchain semantics. CAD-bench
126 is therefore designed around task-specific observables rather than visual complexity alone: visually

Table 1: CAD-bench task taxonomy. Difficulty is authored for benchmark design; as shown later, model difficulty is not strictly monotone in the same order.

Difficulty	Count	Example tasks	Main capability stress
Easy	3	cube, box, ring spacer	pose, units, simple solids
Medium	4	L-bracket, extrusion segment, slotted plate, hex nut blank	feature placement, boolean modeling, conventions
Hard	6	flange, bolted ring, stepped block, large slotted plate, spur gear, M3 socket-head screw	denser geometry, standards-like structure, gears, threads
Functional	4	gearbox, custom threaded pair, compound gearbox, right-angle compound gearbox	assembly reasoning, interface compatibility, rigid-body functionality

127 simple tasks may require exact conventions or mating constraints, while visually complex tasks may
 128 be straightforward once decomposed into deterministic features.

129 4 Benchmark Design and Artifact Construction

130 4.1 Scope and task taxonomy

131 CAD-bench contains 17 tasks divided into four difficulty levels: easy, medium, hard, and functional.
 132 The repository keeps the historical internal label `insane` for the functional bucket, but the paper uses
 133 “functional” because that is the capability being measured. Table 1 summarizes the task structure.

134 The benchmark intentionally mixes single-part and multi-part tasks. The early tasks test whether a
 135 model can synthesize basic parametric geometry. Middle tasks add feature density and standards-like
 136 details. The final tasks require reasoning about interfaces and physical behavior rather than about
 137 isolated solids alone.

138 4.2 Task packaging

139 Each task is stored as a prompt, metadata, expected values, and evaluator configuration. In the
 140 repository, tasks are specified through files such as `prompt.txt`, `task.toml`, and optional reference
 141 assets used for reproducibility. The current public task format includes reference programs in
 142 `gold.py`, but those are benchmark-side reference implementations rather than the definition of the
 143 task itself. The expected metadata includes task ID, difficulty, evaluator, and task-specific expected
 144 values. Public task payloads can be loaded from a local mirror or from a Hugging Face repository.

145 This packaging supports two important use cases: exact reproduction of published harness runs and
 146 evaluation of custom harnesses against the same task set and scoring logic.

147 The tasks are synthetic but engineering-grounded: they specify concrete dimensions, poses, mating
 148 interfaces, thread conventions, and gearbox behavior that a CAD artifact must satisfy. This construc-
 149 tion is intended to make the benchmark sound as an evaluation artifact: the prompts define real CAD
 150 outcomes, the metadata records the expected observables, and the evaluators test submitted artifacts
 151 rather than checking whether a model reproduced a particular text program.

152 4.3 Runtime and submission modes

153 The benchmark runtime supports two submission styles. In one-shot mode, current code-oriented
 154 harnesses return CAD code inside a structured XML block, with the built-in implementation using
 155 `Build123D` [build123d contributors, 2026]. In step-based agent mode, an interactive harness writes
 156 the final CAD artifact to a designated STEP path inside the evaluation environment. In both cases, the
 157 runtime executes the submission in a containerized workspace and converts the result into a common
 158 scoring pipeline.

159 This paper reports complete standalone model runs separately from complete agent runs. The agent
 160 interface remains part of the benchmark runtime, but agent rows are not mixed into the standalone
 161 model table because they use interactive tool feedback and a different wall-clock budget.

162 4.4 Scoring and verification

163 Each task returns normalized metrics including `build_success`, `task_score`, `overall_score`,
164 and a benchmark-facing reward:

$$\text{reward} = 0.05 \cdot \text{build_success} + 0.95 \cdot \text{overall_score}. \quad (1)$$

165 At the benchmark level, CAD-bench first averages task `overall_score` within each difficulty bucket.
166 Missing benchmark tasks are inserted as explicit zero-score rows for result and paper aggregation.
167 The final score is the category-weighted mean

$$S = \frac{1S_{\text{easy}} + 2S_{\text{medium}} + 3S_{\text{hard}} + 4S_{\text{functional}}}{1 + 2 + 3 + 4}. \quad (2)$$

168 This means a zero on a difficulty bucket remains part of the denominator rather than silently dropping
169 that category from the leaderboard. The 1:2:3:4 weighting is a benchmark-design choice rather than
170 a statistical estimate of task importance. It prevents the 13 static tasks from dominating the four
171 functional tasks, while still reporting every bucket separately. We do not use any aggregation rule
172 for statistically significant pairwise model-ranking claims; Section 7 reports a small aggregation-
173 sensitivity check.

174 The task evaluators are designed to match engineering reality. The benchmark uses:

- 175 • mesh- and part-level dimensional checks,
- 176 • pose and placement metrics,
- 177 • reference-geometry gates that downweight large deviations from the authored solution,
- 178 • thread-specific profile and pitch metrics for the custom threaded-pair task, and
- 179 • Blender-based rigid-body simulation for gearbox tasks.

180 Reference-geometry gates are used only as coarse sanity gates on selected tasks. They compare
181 boolean volume differences between the submitted artifact and the reference artifact, then multiply
182 the task score by a smooth penalty when the submitted shape has large missing or extra volume. The
183 task-specific metrics remain responsible for dimensions, poses, threads, and function. This design
184 allows ordinary alternative CAD construction histories, but intentionally reduces credit for artifacts
185 that satisfy a few scalar checks while omitting major geometry.

186 This combination is central to the benchmark. A large part of CAD-bench’s value is that it measures
187 the difference between “the program ran” and “the artifact solved the task.”

188 4.5 Scorer validation and ablations

189 Because CAD-bench uses task-specific evaluators, the release includes explicit scorer-calibration
190 cases in `metadata/cad-bench-scorer-validation.json`. These cases cover reference sub-
191 missions, independently authored valid submissions, minor defects, major defects, and nonbuild
192 submissions across representative easy, hard-static, standards-like, threaded-pair, and functional tasks.
193 Table 2 summarizes the main cases.

194 The validation cases serve two purposes. First, they test invariance to implementation details by
195 giving full credit to independently authored valid artifacts that are not copies of the corresponding
196 `gold.py` programs. Second, they test sensitivity to engineering-relevant defects. Several defective
197 artifacts build successfully but receive low full scores, which demonstrates why build success alone is
198 an inadequate CAD metric. Generic geometry proxies also over-credit local-but-incomplete artifacts,
199 such as a screw without a socket or a flange without bolt holes. For gearbox tasks, the decisive signal
200 is rigid-body behavior: a rigid bridge with plausible axle holes builds and may have the approximately
201 correct bounding box, but receives zero because it cannot transmit the required motion.

202 **Worked M3 example.** The M3 socket-head screw evaluator illustrates the scoring pattern. It first
203 checks pose and dimensions: top of head at $z = 0$, full body in $z \leq 0$, head diameter and height,
204 under-head length, and major diameter. It then checks contact geometry: socket depth, across-flats
205 estimate, open radius at multiple depths, and sixfold harmonic signal. Finally it checks task-specific
206 thread observables: helical signal, pitch, handedness, and approximate turns-to-seat. The final score

Table 2: Representative scorer-calibration cases. Full is the CAD-bench overall_score. Build is a build-success-only ablation. Proxy is the available generic geometry proxy where one exists; gearbox rows are dominated by functional simulation rather than a generic geometry proxy.

Task	Case	Full	Build	Proxy
Cube	independent valid cube	1.000	1.000	1.000
Cube	correct cube shifted off the required pose	0.000	1.000	0.000
Flange	independently constructed valid flange	1.000	1.000	1.000
Flange	bolt circle 1 mm too small	0.800	1.000	0.800
Flange	missing all four bolt holes	0.176	1.000	0.341
M3 screw	independently authored equivalent	1.000	1.000	1.000
M3 screw	threaded screw with missing socket	0.629	1.000	0.723
M3 screw	plain unthreaded shank	0.273	1.000	0.435
Gearbox	reference functional gear pair	0.962	1.000	–
Gearbox	rigid bridge between axle holes	0.000	1.000	–
Gearbox	reference gear pair shifted off axles	0.000	1.000	–
Threaded pair	reference mating threaded pair	0.974	1.000	–
Threaded pair	plausible unthreaded bolt and nut	0.406	1.000	–
Compound gearbox	rigid three-axle bridge	0.000	1.000	–
Right-angle gearbox	intersecting reference-like assembly	0.000	1.000	–

207 is a weighted mean over generic pose, generic dimensions, generic contact, task thread, and task
 208 insertion modules. This is why the independently authored equivalent fixture receives full credit, a
 209 threaded screw without a socket receives partial credit, and an unthreaded cylinder receives only low
 210 residual credit despite building.

211 5 Experimental Protocol

212 5.1 Model and harness matrix

213 The benchmark is intended to support a mixed portfolio of closed and open providers. The current
 214 reviewer result payload contains 19 complete standalone model sweeps and 32 complete agent rows
 215 satisfying the validity criteria in Appendix B.3. Table 3 summarizes these rows in the main paper,
 216 while Appendix B records the exact run identifiers and full row-level results. The current standalone
 217 model set covers OpenAI, DeepSeek, Gemini, Vercel AI Gateway, and free OpenRouter rows.

218 The paper focuses on completed full-benchmark runs, meaning sweeps that cover all 17 tasks
 219 and reach model output plus benchmark scoring. This choice avoids mixing partially completed
 220 campaigns with full evaluations. Provider quota failures, API outages, malformed transport responses,
 221 harness setup failures, and upload failures are retained in the provenance artifacts for accounting but
 222 excluded from model-quality comparisons and from the reviewer result table because they do not
 223 measure the model’s CAD capability.

224 5.2 Reporting protocol

225 For each completed run, we report benchmark score, per-difficulty aggregates, wall-clock time, and
 226 estimated cost when the harness exposes enough usage and pricing data. When multiple complete
 227 runs exist for the same provider, model, access mode, and reasoning level, all complete rows in the
 228 reviewer result payload are reported, and the exact run identifiers are recorded in Appendix B. The
 229 current paper does not claim statistically significant pairwise rankings.

230 This choice reflects the current state of the experimental log rather than an intended limitation of the
 231 benchmark. CAD-bench itself stores the artifacts needed for repeated-run reporting, and larger future
 232 campaigns can replace unreplicated full-sweep rows with repeated-run aggregates.

233 6 Results: What CAD-bench Reveals

234 6.1 Overall performance

235 Table 3 summarizes the complete full-benchmark evaluations mirrored for anonymous review. Full
 236 leaderboards, run identifiers, wall-clock times, and cost estimates are reported in Appendix B; the

Table 3: Summary of complete CAD-bench baselines. Scores are category-weighted using Eq. (2). Full per-row results and run identifiers are reported in Appendix B.

Setting	Rows	Best overall	Mean easy	Mean hard	Mean functional
Standalone models	19	0.599	0.930	0.585	0.037
Agent harnesses	32	0.677	0.906	0.761	0.141

237 main text focuses on aggregate conclusions because many configurations currently have only one
238 complete sweep.

239 Three conclusions emerge from the completed rows.

240 First, CAD-bench is not saturated. The best current standalone model reaches 0.599 overall, and the
241 best agent harness reaches 0.677. These scores leave substantial headroom despite high performance
242 on many easy and static tasks.

243 Second, static geometry is substantially easier than functional CAD. Completed standalone runs
244 average 0.585 on hard-static tasks but only 0.037 on functional tasks. Agent harnesses improve
245 both categories, but the same pattern remains: they average 0.761 on hard-static tasks and 0.141 on
246 functional tasks.

247 Third, the benchmark distinguishes model and harness settings. Agent harnesses generally outperform
248 standalone one-shot generation, consistent with the hypothesis that execution feedback helps repair
249 dimensional, geometric, and interface errors. However, agent success is uneven, and functional
250 assemblies remain difficult even when submissions build successfully.

251 The difficulty claim is not based only on the names of the tiers. Under the same reported artifacts,
252 replacing CAD-bench’s full score with weaker proxies would make the benchmark look easier.
253 Reported agent submissions build the right-angle gearbox in 96.9% of attempts but average only
254 4.0% full score, and standalone submissions build the gearbox task in 63.2% of attempts but average
255 only 6.7% full score. Thus the functional tasks are not merely harder by our label, they are clearly
256 also harder for the language models.

257 6.2 Per-difficulty behavior

258 Across the 19 reported standalone model sweeps, the mean score by difficulty is 0.930 on easy,
259 0.753 on medium, 0.585 on hard, and 0.037 on functional tasks. Across the 32 reported agent
260 sweeps, the corresponding means are 0.906, 0.849, 0.761, and 0.141. The drop on functional tasks
261 is expected. The remaining drop from medium to hard indicates that standards-like details, dense
262 feature placement, and gear/thread structure still matter after simple part synthesis is solved.

263 In CAD-bench, difficulty labels are benchmark-design labels, not claims about a strict monotone
264 model ordering. Several hard tasks are still static, deterministic solids with many dimensions but
265 little ambiguity, while some lower-level tasks are sensitive to conventions such as exact orientation,
266 alignment, or hole placement. The current results therefore reinforce the point that CAD evaluation
267 must be task-specific rather than purely complexity-labeled.

268 7 Discussion and Limitations

269 CAD-bench is mainly an evaluation artifact rather than a final model leaderboard. Its main claim
270 is that executable CAD scoring exposes failures that build success, rendering validity, and coarse
271 geometry proxies miss. The reported rows are useful diagnostic baselines, but many configurations
272 have only one full sweep. We therefore avoid statistically significant pairwise model-ranking claims.
273 Repeated-run statistics in Appendix B.2 support the static-versus-functional difficulty contrast for
274 seven standalone configurations, and future leaderboard versions should promote rows only after
275 repeated full sweeps under a declared aggregation rule.

276 The aggregation rule is not significantly driving the main benchmark either. Under the current
277 complete result payload, the top standalone row is unchanged under the paper’s weighted score, equal
278 bucket averaging, and task-count weighting: Gemini 3.1 Pro offline scores 0.599, 0.727, and 0.709
279 respectively. For agent rows, weighted and equal-bucket aggregation both put Codex GPT-5.4 mini

Table 4: Per-task mean overall_score and build success across the complete reported rows. Standalone means use 19 model sweeps; agent means use 32 agent sweeps. This table is computed from the 867 task rows in the full reviewer artifact, not from the difficulty aggregates alone.

Difficulty	Task	Standalone score	Standalone build	Agent score	Agent build
Easy	cube	0.895	0.947	0.938	0.969
Easy	box	0.947	1.000	0.844	0.844
Easy	ring spacer	0.947	1.000	0.938	0.938
Medium	L-bracket	0.632	0.842	0.906	0.969
Medium	extrusion	0.868	0.947	0.875	0.875
Medium	slotted plate	0.672	0.947	0.864	0.969
Medium	hex nut blank	0.842	0.895	0.750	0.906
Hard	flange	0.746	0.842	0.906	0.969
Hard	bolted ring	0.789	0.789	0.947	1.000
Hard	stepped block	0.812	0.947	1.000	1.000
Hard	large slotted plate	0.718	0.842	0.853	1.000
Hard	spur gear	0.237	0.526	0.345	0.844
Hard	M3 socket screw	0.207	0.526	0.512	0.875
Functional	gearbox	0.067	0.632	0.316	0.844
Functional	threaded pair	0.012	0.211	0.071	0.938
Functional	compound gearbox	0.069	0.632	0.138	0.719
Functional	right-angle gearbox	0.000	0.158	0.040	0.969

280 medium web first, at 0.677 and 0.781; task-count weighting instead puts Pi GPT-5.4 high offline first
 281 at 0.772. The qualitative conclusion is stable: easy/static CAD is much stronger than functional CAD,
 282 while fine-grained row ordering should remain diagnostic until repeated runs.

283 The benchmark is also small. Seventeen tasks are enough for a high-information initial release
 284 that covers primitives, feature placement, standards-like details, threaded interfaces, and simulated
 285 mechanisms, but they are not enough to claim broad CAD-agent competence. The public release is
 286 therefore intended to be versioned. Appendix A describes the planned lifecycle: immutable public
 287 task bundles, private refresh tasks with canary hashes, contamination monitoring, and repeated-run
 288 promotion for leaderboard rows.

289 Task-specific scorers encode engineering assumptions. This is necessary for CAD: a threaded pair,
 290 socket-head screw, or gearbox cannot be judged by build success alone. It also creates a risk that a
 291 valid alternative solution could miss an authored observable. The scorer-validation cases in Table 2
 292 partially address this by giving full credit to independently authored valid artifacts and low credit
 293 to controlled defects, but they are not exhaustive. The intended use of CAD-bench is therefore
 294 comparative and diagnostic: inspect per-task traces, scorer components, and validation cases, not
 295 only the aggregate score.

296 8 Conclusion

297 CAD-bench is an executable benchmark for CAD generation that emphasizes artifact correctness over
 298 surface plausibility. The current v0 release shows that models and agent harnesses can often produce
 299 runnable CAD, but still fail on standards-like details, mating interfaces, and functional assemblies.
 300 The benchmark’s main contribution is the scoring interface and reproducible artifact protocol: it
 301 separates easy geometry from mechanically meaningful CAD behavior and gives future systems a
 302 concrete target for improving assemblies and interfaces that actually work.

303 References

- 304 Kamel Alrashedy, Pradyumna Tambwekar, Zulfiqar Zaidi, Megan Langwasser, Wei Xu, and
305 Matthew Gombolay. Generating CAD Code with Vision-Language Models for 3D Designs.
306 arXiv:2410.05340, 2024. URL: <https://arxiv.org/abs/2410.05340>.
- 307 Jesse Barkley, Rumi Loghmani, and Amir Barati Farimani. CADSmith: Multi-Agent CAD Generation
308 with Programmatic Geometric Validation. arXiv:2603.26512, 2026. URL: <https://arxiv.org/abs/2603.26512>.
- 310 build123d contributors. build123d documentation. 2026. URL: <https://build123d.readthedocs.io/en/stable/index.html>.
- 312 CAD Arena. CAD Arena: Open Benchmark for AI-Generated Parametric CAD. 2026. URL:
313 <https://cadarena.dev/>.
- 314 Jun Shern Chan, Neil Chowdhury, Oliver Jaffe, James Aung, Dane Sherburn, Evan Mays, Giulio
315 Starace, Kevin Liu, Leon Maksin, Tejal Patwardhan, Lilian Weng, and Aleksander Mądry. MLE-
316 bench: Evaluating Machine Learning Agents on Machine Learning Engineering. arXiv:2410.07095,
317 2024. URL: <https://arxiv.org/abs/2410.07095>.
- 318 Anna C. Doris, Daniele Grandi, Ryan Tomich, Md Ferdous Alam, Mohammadmehdi Ataei, Hyunmin
319 Cheong, and Faez Ahmed. DesignQA: A Multimodal Benchmark for Evaluating Large Language
320 Models' Understanding of Engineering Documentation. arXiv:2404.07917, 2024. URL: <https://arxiv.org/abs/2404.07917>.
- 322 Epoch AI. CadEval. 2026. URL: <https://epoch.ai/benchmarks/cad-eval>.
- 323 Yaroslav Ganin, Sergey Bartunov, Yujia Li, Ethan Keller, and Stefano Saliceti. Computer-Aided
324 Design as Language. arXiv:2105.02769, 2021. URL: <https://arxiv.org/abs/2105.02769>.
- 325 Yandong Guan, Xilin Wang, Xingxi Ming, Jing Zhang, Dong Xu, and Qian Yu. CAD-Coder: Text-to-
326 CAD Generation with Chain-of-Thought and Geometric Reward. arXiv:2505.19713, 2025. URL:
327 <https://arxiv.org/abs/2505.19713>.
- 328 Mohammad Sadil Khan, Sankalp Sinha, Talha Uddin Sheikh, Didier Stricker, Sk Aziz Ali, and
329 Muhammad Zeshan Afzal. Text2CAD: Generating Sequential CAD Models from Beginner-to-
330 Expert Level Text Prompts. In *Advances in Neural Information Processing Systems*, volume 37,
331 2024. URL: <https://openreview.net/forum?id=5k9XeHIK3L>.
- 332 Carlos E. Jimenez, John Yang, Alexander Wettig, Shunyu Yao, Kexin Pei, Ofir Press, and
333 Karthik Narasimhan. SWE-bench: Can Language Models Resolve Real-World GitHub Issues?
334 arXiv:2310.06770, 2023. URL: <https://arxiv.org/abs/2310.06770>.
- 335 Jing Yu Koh, Robert Lo, Lawrence Jang, Vikram Duvvur, Ming Chong Lim, Po-Yu Huang, Graham
336 Neubig, Shuyan Zhou, Ruslan Salakhutdinov, and Daniel Fried. VisualWebArena: Evaluating
337 Multimodal Agents on Realistic Visual Web Tasks. arXiv:2401.13649, 2024. URL: <https://arxiv.org/abs/2401.13649>.
- 339 Xiao Liu, Hao Yu, Hanchen Zhang, Yifan Xu, Xuanyu Lei, Hanyu Lai, Yu Gu, Hangliang Ding,
340 Kaiwen Men, Kejuan Yang, Shudan Zhang, Xiang Deng, Aohan Zeng, Zhengxiao Du, Chenhui
341 Zhang, Sheng Shen, Tianjun Zhang, Yu Su, Huan Sun, Minlie Huang, Yuxiao Dong, and Jie
342 Tang. AgentBench: Evaluating LLMs as Agents. arXiv:2308.03688, 2023. URL: <https://arxiv.org/abs/2308.03688>.
- 344 Gregoire Mialon, Clementine Fourier, Craig Swift, Thomas Wolf, Yann LeCun, and Thomas
345 Scialom. GAIA: a benchmark for General AI Assistants. arXiv:2311.12983, 2023. URL:
346 <https://arxiv.org/abs/2311.12983>.
- 347 Anka Reuel, Amelia Hardy, Chandler Smith, Max Lamparth, Malcolm Hardy, and Mykel J. Kochen-
348 derfer. BetterBench: Assessing AI Benchmarks, Uncovering Issues, and Establishing Best Practices.
349 arXiv:2411.12990, 2024. URL: <https://arxiv.org/abs/2411.12990>.

- 350 Ari Seff, Yaniv Ovadia, Wenda Zhou, and Ryan P. Adams. SketchGraphs: A Large-Scale Dataset
351 for Modeling Relational Geometry in Computer-Aided Design. arXiv:2007.08506, 2020. URL:
352 <https://arxiv.org/abs/2007.08506>.
- 353 Giulio Starace, Oliver Jaffe, Dane Sherburn, James Aung, Jun Shern Chan, Leon Maksin, Rachel Dias,
354 Evan Mays, Benjamin Kinsella, Wyatt Thompson, Johannes Heidecke, Amelia Glaese, and Tejal
355 Patwardhan. PaperBench: Evaluating AI’s Ability to Replicate AI Research. arXiv:2504.01848,
356 2025. URL: <https://arxiv.org/abs/2504.01848>.
- 357 Colin White, Samuel Dooley, Manley Roberts, Arka Pal, Ben Feuer, Siddhartha Jain, Ravid Shwartz-
358 Ziv, Neel Jain, Khalid Saifullah, Siddhartha Naidu, Chinmay Hegde, Yann LeCun, Tom Goldstein,
359 Willie Neiswanger, and Micah Goldblum. LiveBench: A Challenging, Contamination-Free LLM
360 Benchmark. arXiv:2406.19314, 2024. URL: <https://arxiv.org/abs/2406.19314>.
- 361 Haoyang Xie and Feng Ju. Text-to-CadQuery: A New Paradigm for CAD Generation with Scalable
362 Large Model Capabilities. arXiv:2505.06507, 2025. URL: <https://arxiv.org/abs/2505.06507>.
- 363
- 364 Karl D. D. Willis, Yewen Pu, Jieliang Luo, Hang Chu, Tao Du, Joseph G. Lambourne, Armando
365 Solar-Lezama, and Wojciech Matusik. Fusion 360 Gallery: A Dataset and Environment for
366 Programmatic CAD Construction from Human Design Sequences. *ACM Transactions on Graphics*,
367 40(4), 2021. URL: <https://arxiv.org/abs/2010.02392>.
- 368 Hjalmar Wijk, Tao Lin, Joel Becker, Sami Jawhar, Neev Parikh, Thomas Broadley, Lawrence Chan,
369 Michael Chen, Josh Clymer, Jai Dhyani, Elena Elicheva, Katharyn Garcia, Brian Goodrich,
370 Nikola Jurkovic, Megan Kinniment, Aron Lajko, Seraphina Nix, Lucas Sato, William Saunders,
371 Maksym Taran, Ben West, and Elizabeth Barnes. RE-Bench: Evaluating frontier AI R&D
372 capabilities of language model agents against human experts. arXiv:2411.15114, 2024. URL:
373 <https://arxiv.org/abs/2411.15114>.
- 374 Rundi Wu, Chang Xiao, and Changxi Zheng. DeepCAD: A Deep Generative Network for Computer-
375 Aided Design Models. In *Proceedings of the IEEE/CVF International Conference on Computer*
376 *Vision*, 2021. URL: <https://arxiv.org/abs/2105.09492>.
- 377 Tianbao Xie, Danyang Zhang, Jixuan Chen, Xiaochuan Li, Siheng Zhao, Ruisheng Cao, Toh Jing
378 Hua, Zhoujun Cheng, Dongchan Shin, Fangyu Lei, Yitao Liu, Yiheng Xu, Shuyan Zhou, Silvio
379 Savarese, Caiming Xiong, Victor Zhong, and Tao Yu. OSWorld: Benchmarking Multimodal
380 Agents for Open-Ended Tasks in Real Computer Environments. arXiv:2404.07972, 2024. URL:
381 <https://arxiv.org/abs/2404.07972>.

382 A Benchmark lifecycle, release policy, and broader impacts

383 The artifact is designed to support a versioned benchmark lifecycle rather than a single frozen
384 leaderboard. The submitted release should be interpreted as CAD-bench v0: a public, reproducible
385 task set with exact run artifacts. Future leaderboard versions should use four controls.

- 386 • **Frozen public splits.** Published task bundles are immutable and identified by manifest
387 hashes and dataset revisions.
- 388 • **Private refresh tasks.** New functional and standards-like tasks can be staged privately, with
389 salted bundle-hash canaries exported before evaluation to establish that the tasks existed
390 before result collection.
- 391 • **Versioned leaderboards.** Rows should be reported against an explicit task-set version; old
392 rows should not be silently mixed with refreshed tasks.
- 393 • **Repeated-run promotion.** A configuration should move from diagnostic reporting to a
394 primary leaderboard row only after independent full sweeps are available with a declared
395 aggregation rule.

396 This policy is already reflected in the release files: task manifests record per-file hashes, the release
397 pins the task dataset revision, and the source artifact contains a canary-export command for private

398 or unreleased tasks. The current paper still reports v0 diagnostic rows because the full repeated-run
 399 campaign is not complete; this is why the claims are framed around benchmark behavior and failure
 400 modes rather than statistically significant pairwise model rankings.

401 CAD agents could provide several benefits: lowering the barrier to prototyping, supporting
 402 accessibility-focused device customization, improving engineering education, and reducing repetitive
 403 modeling work in low-volume manufacturing workflows.

404 The same capabilities also create risks. More capable CAD generation could accelerate the design of
 405 unsafe hardware, restricted components, or deceptive engineering artifacts. It may also encourage
 406 overtrust in automatically generated mechanical designs that have not received appropriate engi-
 407 neering review. Transparent, function-aware benchmarking is one mitigation: it makes it harder to
 408 confuse runnable CAD code with mechanically correct design.

409 B Reproducibility and Reporting Details

410 B.1 Raw complete result rows

411 Tables 5–7 record the exact complete rows summarized in Table 3. Scores are reported as percentages
 412 for readability. The run column is the run identifier corresponding to the stored report artifact. Partial
 413 runs and infrastructure-failed rows are excluded from model-quality comparisons.

Table 5: Raw complete standalone model rows.

Rank	Harness	Overall	Easy	Med.	Hard	Func.	Min.	Cost	Run
1	gemini/gemini-3.1-pro-preview-offline-none	59.9	100.0	96.8	70.8	23.2	25.4	5.8625	gemini__gemini-3.1-pro-preview-offline-none_20260428T194054Z_2313253
2	openai/gpt-5.4-pro-web-med	53.9	100.0	96.8	78.5	2.5	325.737.0405	openai__gpt-5.4-pro-web-med_20260423T122643Z_1626647	
3	vercel/alibaba/qwen-3.6-max-preview-offline-none	52.8	100.0	96.8	76.2	1.5	42.8	3.2293	vercel__alibaba_qwen-3.6-max-preview-offline-none_20260430T003706Z_2473332
4	gemini/gemini-3-flash-preview-offline-none	52.6	100.0	96.8	77.4	0.0	21.1	1.3596	gemini__gemini-3-flash-preview-offline-none_20260428T121645Z_2280737
5	openrouter/tencent/hy3-preview:free-offline-none	51.1	100.0	96.8	60.8	8.7	51.4	-	openrouter__tencent_hy3-preview-free-offline-none_20260428T035110Z_2240039
6	openai/gpt-5.3-codex-spark-offline-med	50.4	100.0	96.8	70.0	0.0	6.1	4.9685	openai__gpt-5.3-codex-spark-offline-med_20260427T023421Z_2115684
7	vercel/alibaba/qwen3.6-27b-offline-none	50.0	100.0	96.8	68.9	0.0	58.6	1.5128	vercel__alibaba_qwen3.6-27b-offline-none_20260430T003707Z_2473345
8	openai/gpt-5.3-codex-spark-offline-xhigh	48.9	100.0	96.8	65.0	0.0	6.0	5.4877	openai__gpt-5.3-codex-spark-offline-xhigh_20260427T024838Z_2122296
9	openai/gpt-5.5-offline-med	48.7	100.0	71.8	65.0	12.0	9.8	8.6697	openai__codex__gpt-5.5-offline-med_20260426T055852Z_2010916
10	openai/gpt-5.5-offline-low	47.5	100.0	71.8	76.3	0.6	10.1	9.7384	openai__codex__gpt-5.5-offline-low_20260426T054800Z_2000292
11	deepseek/deepseek-v4-pro-offline-none	47.5	100.0	59.4	68.0	12.9	55.7	3.6274	deepseek__deepseek-v4-pro-offline-none_20260426T042645Z_1958729
12	openai/gpt-5.5-offline-high	45.6	66.7	71.8	82.1	0.0	8.9	7.8075	openai__codex__gpt-5.5-offline-high_20260426T061252Z_2019625
13	openai/gpt-5.5-offline-xhigh	45.6	100.0	71.8	59.3	8.7	9.9	9.2523	openai__codex__gpt-5.5-offline-xhigh_20260426T062226Z_2023449
14	openai/gpt-5.3-codex-spark-offline-high	40.4	100.0	96.8	36.9	0.0	6.1	5.2469	openai__gpt-5.3-codex-spark-offline-high_20260427T024046Z_2118961
15	deepseek/deepseek-v4-flash-offline-none	39.5	100.0	64.3	55.5	0.0	18.1	0.2778	deepseek__deepseek-v4-flash-offline-none_20260426T052907Z_1985868
16	gemini/gemini-3.1-flash-lite-preview-offline-none	34.2	100.0	75.0	30.8	0.0	19.2	0.1993	gemini__gemini-3.1-flash-lite-preview-offline-none_20260428T050032Z_2249465
17	openai/gpt-5.3-codex-spark-offline-low	33.9	100.0	71.8	31.7	0.0	5.8	4.2688	openai__gpt-5.3-codex-spark-offline-low_20260427T022715Z_2111735
18	openrouter/inclusionai/ling-2.6-1t:free-offline-none	13.9	33.3	2.7	33.3	0.0	12.8	-	openrouter__inclusionai_ling-2.6-1t-free-offline-none_20260428T033043Z_2234914
19	gemini/gemma-3-27b-it-offline_nodocs-none	8.1	66.7	0.0	4.6	0.0	3.7	-	gemini__gemma-3-27b-it-offline_nodocs-none_20260430T022138Z_2493048

414 B.2 Repeated-run statistical check

415 The current paper reports complete full 17-task sweeps for every row, but many configurations
 416 have only one independent sweep. The paper therefore does not use those rows for pairwise
 417 model-ranking significance claims. To test the main qualitative difficulty claim under repeated
 418 sampling, the source artifact includes 47 complete local reports, with repeated full-sweep statistics
 419 for seven standalone configurations, summarized in metadata/cad-bench-repeat-stats.json.
 420 Additional attempted repeats across the reported standalone and agent matrix were excluded when

Table 6: Raw complete agent rows.

Rank	Agent Harness	Overall	Easy	Med.	Hard	Func.	Min.	Cost	Run
1	Codex codex/gpt-5.4-mini-web-med	67.7	100.0	96.8	78.7	36.8	234.0	2.4446	codex_gpt-5.4-mini-web-med_20260420T000623Z_676664
2	Pi pi/gpt-5.4-offline-high	67.1	100.0	96.8	82.9	32.1	618.8	4.6416	pi_gpt-5.4-offline-high_20260422T092300Z_1340452
3	Codex codex/gpt-5.4-offline-high	65.3	66.7	71.8	83.6	47.9	245.2	8.2311	codex_gpt-5.4-offline-high_20260419T190624Z_603644
4	Pi pi/gpt-5.4-offline-xhigh	64.7	100.0	96.8	88.0	22.3	418.3	4.8298	pi_gpt-5.4-offline-xhigh_20260421T075050Z_1041869
5	Pi pi/gpt-5.4-web-high	64.6	100.0	96.8	82.9	25.9	281.6	4.5650	pi_gpt-5.4-web-high_20260422T092300Z_1340450
6	Pi pi/gpt-5.4-mini-web-xhigh	61.0	100.0	100.0	84.5	14.2	240.2	2.4442	pi_gpt-5.4-mini-web-xhigh_20260422T092301Z_1340454
7	Pi pi/gpt-5.4-mini-offline-high	60.3	100.0	75.0	80.3	28.0	169.0	3.5843	pi_gpt-5.4-mini-offline-high_20260422T151655Z_1441926
8	Pi pi/gpt-5.4-offline-med	59.9	100.0	96.8	76.2	19.3	304.2	1.9074	pi_gpt-5.4-offline-med_20260421T075052Z_1041872
9	Codex codex/gpt-5.4-web-low	59.1	100.0	96.8	79.4	14.8	279.8	6.3453	codex_gpt-5.4-web-low_20260419T190502Z_603630
10	Codex codex/gpt-5.4-web-med	59.1	100.0	71.8	78.0	28.2	260.7	5.7144	codex_gpt-5.4-web-med_20260419T190540Z_603629
11	Pi pi/gpt-5.4-mini-offline-xhigh	59.0	100.0	96.8	83.5	11.4	282.4	3.4776	pi_gpt-5.4-mini-offline-xhigh_20260422T121627Z_1402824
12	Codex codex/gpt-5.4-mini-offline-med	57.8	100.0	96.8	52.5	31.7	157.2	3.3316	codex_gpt-5.4-mini-offline-med_20260420T004107Z_683812
13	Codex codex/gpt-5.4-web-xhigh	56.9	100.0	96.8	89.7	1.7	396.2	12.7137	codex_gpt-5.4-web-xhigh_20260421T075051Z_1041874
14	Codex codex/gpt-5.4-mini-offline-xhigh	56.9	100.0	96.8	77.8	10.5	120.9	4.0830	codex_gpt-5.4-mini-offline-xhigh_20260422T205507Z_1497871
15	Pi pi/gpt-5.4-web-xhigh	56.0	100.0	71.8	82.9	16.9	419.3	4.7572	pi_gpt-5.4-web-xhigh_20260421T075052Z_1041870
16	Pi pi/gpt-5.4-mini-web-med	54.5	100.0	96.8	81.8	1.5	191.2	1.6381	pi_gpt-5.4-mini-web-med_20260422T092300Z_1340457

Table 7: Raw complete agent rows, continued.

Rank	Agent Harness	Overall	Easy	Med.	Hard	Func.	Min.	Cost	Run
17	Pi pi/gpt-5.4-mini-web-high	54.0	100.0	96.8	75.7	4.7	130.5	2.5139	pi_gpt-5.4-mini-web-high_20260422T18318Z_1468372
18	Codex codex/gpt-5.4-mini-web-high	53.8	100.0	96.8	78.3	2.3	389.9	3.6274	codex_gpt-5.4-mini-web-high_20260420T235625Z_845623
19	Codex codex/gpt-5.4-offline-xhigh	53.5	33.3	71.8	76.2	32.3	289.6	11.4765	codex_gpt-5.4-offline-xhigh_20260419T190535Z_603631
20	Pi pi/gpt-5.4-offline-low	53.3	100.0	75.0	84.4	7.6	171.2	1.5744	pi_gpt-5.4-offline-low_20260419T061026Z_383966
21	Pi pi/gpt-5.4-web-med	53.0	100.0	71.8	80.1	11.6	366.2	2.1871	pi_gpt-5.4-web-med_20260421T075053Z_1041868
22	Pi pi/gpt-5.4-web-low	52.6	100.0	96.8	76.2	0.9	65.1	1.1919	pi_gpt-5.4-web-low_20260422T205357Z_1497569
23	Codex codex/gpt-5.4-web-high	52.1	66.7	96.8	80.1	5.0	300.5	7.9083	codex_gpt-5.4-web-high_20260419T190546Z_603612
24	Codex codex/gpt-5.4-mini-offline-high	51.3	100.0	96.8	71.1	1.5	169.1	4.4820	codex_gpt-5.4-mini-offline-high_20260420T004018Z_683339
25	Codex codex/gpt-5.4-mini-web-xhigh	51.1	100.0	71.8	86.9	1.6	416.0	5.9764	codex_gpt-5.4-mini-web-xhigh_20260421T075051Z_1041875
26	Pi pi/gpt-5.4-mini-offline-med	49.7	100.0	96.8	64.5	2.4	241.8	2.4102	pi_gpt-5.4-mini-offline-med_20260422T094713Z_1351194
27	Pi pi/gpt-5.4-mini-web-low	48.1	100.0	75.0	77.1	0.0	189.4	0.5317	pi_gpt-5.4-mini-web-low_20260419T064631Z_413438
28	Codex codex/gpt-5.4-mini-offline-low	48.1	100.0	100.0	57.5	2.2	111.3	1.6688	codex_gpt-5.4-mini-offline-low_20260420T014421Z_709021
29	Codex codex/gpt-5.4-offline-low	46.8	33.3	46.8	88.4	18.9	181.1	4.4749	codex_gpt-5.4-offline-low_20260419T190806Z_603646
30	Pi pi/gpt-5.4-mini-offline-low	37.9	100.0	71.8	38.5	4.9	132.1	0.5499	pi_gpt-5.4-mini-offline-low_20260422T100629Z_1359611
31	Codex codex/gpt-5.4-offline-med	36.2	33.3	46.8	61.6	12.6	260.4	8.0969	codex_gpt-5.4-offline-med_20260419T190604Z_603645
32	Codex codex/gpt-5.4-mini-web-low	32.9	66.7	50.0	54.1	0.0	246.5	1.8227	codex_gpt-5.4-mini-web-low_20260420T000933Z_677450

421 provider quota, authentication, unavailable-model, or timeout failures prevented a complete 17-task
 422 report.

423 Across six openai/gpt-5.3-codex-spark-offline-med sweeps, the benchmark score is 0.471
 424 with bootstrap 95% CI [0.453, 0.490]. The easy bucket is stable at 1.000, medium is 0.973 with CI
 425 [0.968, 0.984], hard is 0.572 with CI [0.507, 0.643], and functional is 0.011 with CI [0.000, 0.033].
 426 A paired exact two-sided sign-flip test over full-sweep difficulty aggregates gives $p = 0.03125$ for
 427 hard score minus functional score, with mean difference 0.561 and bootstrap 95% CI [0.486, 0.641].

428 Across the four six-sweep `openai/gpt-5.5-offline-*` repeated checks, benchmark means are
429 0.477 for low with CI [0.455, 0.492], 0.444 for medium with CI [0.391, 0.503], 0.505 for high with
430 CI [0.454, 0.563], and 0.573 for xhigh with CI [0.528, 0.612]. Their hard-minus-functional mean
431 differences are 0.701, 0.483, 0.695, and 0.694, respectively, and each exact sign-flip test gives
432 $p = 0.03125$. The xhigh row has the strongest repeated GPT-5.5 mean here, but these runs still
433 support only a within-configuration static-versus-functional difficulty contrast, not a statistically
434 significant model-ranking claim.

435 Across six `openrouter/tencent/hy3-preview:free-offline-none` sweeps, the benchmark
436 score is 0.518 with CI [0.471, 0.557]; easy is 1.000, medium is 0.823 with CI [0.729, 0.917], hard
437 is 0.664 with CI [0.585, 0.737], and functional is 0.135 with CI [0.095, 0.176]. The hard-minus-
438 functional sign-flip test gives $p = 0.03125$, mean difference 0.529, and CI [0.472, 0.594]. Across five
439 `openrouter/inclusionai/ling-2.6-1t:free-offline-none` sweeps, the benchmark score is
440 0.154 with CI [0.080, 0.236]; hard-minus-functional remains positive at 0.227 with CI [0.084, 0.407],
441 but the exact sign-flip test is $p = 0.125$ at $n = 5$. These repeated-run checks support the paper’s main
442 diagnostic claim that functional CAD remains much harder than static CAD for multiple standalone
443 harnesses, while leaving the model-level ranking tables as diagnostic baselines rather than statistically
444 significant pairwise comparisons.

445 The source artifact also retains the earlier `metadata/cad-bench-repeat-check.json` two-run
446 sanity record for traceability. Future leaderboard versions should report repeated sweeps in the
447 following format:

- 448 • number of independent full sweeps per configuration,
- 449 • central estimate (mean or median),
- 450 • uncertainty interval (e.g., standard deviation or bootstrap 95% CI),
- 451 • policy for selecting primary leaderboard rows when multiple repetitions exist.

452 The current analysis sections avoid pairwise significance claims and use the complete sweeps as
453 baseline diagnostics.

454 **B.3 Run validity criteria**

455 A run is eligible for a complete-results table if it satisfies all of the following conditions:

- 456 • it uses a standalone model harness or an agent harness,
- 457 • it attempts the full 17-task benchmark under one fixed provider, model, access mode, and
458 reasoning setting,
- 459 • each task reaches model-output capture, STEP export when applicable, and benchmark
460 scoring, and
- 461 • failures, if any, are model-output or scored-artifact failures rather than provider quota, API
462 outage, harness setup, or artifact-upload failures.

463 Provider and infrastructure failures are retained in the experiment logs for auditability, but they are not
464 assigned model-quality scores because they do not evaluate the submitted model on CAD generation.

465 **B.4 Compute and environment reporting**

466 The benchmark runtime uses Dockerized environments, Python 3.11, NumPy, trimesh, and Blender
467 for the functional simulation tasks, together with Build123D in the current code-oriented reference
468 and harness pipeline. The reported local runs used the following environment:

- 469 • host hardware: Intel Core i5-7300U CPU, 15 GiB RAM, CPU-only scoring except for
470 Blender’s local simulation/rendering stack,
- 471 • software versions: Docker 29.4.0, Blender 5.1.1, uv 0.11.7, project Python 3.11 environment
472 from `uv.lock`,
- 473 • standalone model wall-clock range in the reported rows: 3.7–325.7 minutes,
- 474 • agent wall-clock range in the reported rows: 65.1–618.8 minutes,

475 • failed-run accounting: provider quota failures are kept in logs but excluded from model-
476 quality comparisons.

477 **B.5 Code, data, and asset availability**

478 The benchmark runtime, scoring code, and harness interface are available in the anonymous source
479 artifact accompanying this paper. Public task payloads are hosted in an anonymous reviewer data
480 mirror, and the reviewer mirror includes a compact result artifact with every complete run reported in
481 this paper.

482 For reviewer inspection, the artifact has a no-API smoke path: `uv run pytest`
483 `tests/test_anonymous_artifacts.py tests/test_scorer_validation_artifact.py`
484 checks the anonymous packaging and scorer-calibration artifact, and `uv run python`
485 `scripts/export_scorer_validation.py` regenerates the calibration JSON. Full model
486 sweeps require provider credentials, but scorer validation and artifact consistency do not.

487 The anonymous submission release pointers are:

- 488 • anonymized source archive: included in the supplemental material and mirrored with the
489 anonymous dataset artifact,
- 490 • anonymized public task dataset mirror: included in the supplemental material,
- 491 • task dataset revision used for this paper: `20ad0b586b70ac2c5e6fe2b501386cf26a887137`,
- 492 • reviewer result artifact: `results/cad-bench-reported-results.json` in the anony-
493 mous dataset mirror,
- 494 • scorer validation artifact: `results/cad-bench-scorer-validation.json` in the anony-
495 mous dataset mirror,
- 496 • run artifact bundles: `provenance/<run-id>/report.json` entries in the anonymous
497 reviewer artifact,
- 498 • provenance URLs: redacted from submitted metadata when they would reveal non-
499 anonymous infrastructure.

500 The full reviewer artifact includes Croissant metadata with Responsible AI fields describing intended
501 use, limitations, task authoring, and safety-relevant non-use cases. This matches the 2026 Evaluations
502 & Datasets track expectation that datasets document how they support evaluative claims rather than
503 serving only as raw files.

504 **B.6 Licenses and upstream assets**

505 Table 8 documents the main upstream software and data assets used by the benchmark runtime and
506 reviewer dataset.

Table 8: Primary upstream software assets used by CAD-bench.

Asset	License / terms	Role in benchmark
Build123D	Apache-2.0	current reference implementation and code-harness dependency, not the benchmark definition itself
Blender	GNU GPL	rigid-body simulation and rendering for functional gearbox evaluation
NumPy	BSD-style license	numeric processing in evaluators and utilities
trimesh	MIT	mesh conversion and geometric analysis
Repository code	MIT	benchmark runtime, harnesses, and scoring stack
Public task dataset	MIT	hosted task prompts, metadata, reference programs, and authored benchmark assets
Authored CAD fixtures	MIT	synthetic task-side fixtures, including the M3 socket-head screw equivalent fixture
Closed model APIs	provider terms	external inference services used for benchmarked models

507 The benchmark tasks and authored fixtures are synthetic benchmark assets rather than redistributed
508 proprietary CAD files.

509 **NeurIPS Paper Checklist**

510 **1. Claims**

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

513 Answer: [Yes]

514 Justification: The abstract and Sections 1 and 4 state the paper’s scope as introducing an
515 executable CAD benchmark, documenting its scoring stack, and reporting the currently com-
516 pleted baseline results with explicit limitations around unreplicated full-sweep comparisons.

517 Guidelines:

- 518 • The answer [N/A] means that the abstract and introduction do not include the claims
519 made in the paper.
- 520 • The abstract and/or introduction should clearly state the claims made, including the
521 contributions made in the paper and important assumptions and limitations. A [No] or
522 [N/A] answer to this question will not be perceived well by the reviewers.
- 523 • The claims made should match theoretical and experimental results, and reflect how
524 much the results can be expected to generalize to other settings.
- 525 • It is fine to include aspirational goals as motivation as long as it is clear that these goals
526 are not attained by the paper.

527 2. Limitations

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

529 Answer: [Yes]

530 Justification: The main results section and Appendix A discuss benchmark size, the current
531 emphasis of the reported baselines on Build123D-based harnesses rather than the benchmark
532 definition itself, incomplete all-provider campaign coverage, current dependence on a small
533 number of completed full sweeps, benchmark exposure risk, and the planned refresh/holdout
534 protocol.

535 Guidelines:

- 536 • The answer [N/A] means that the paper has no limitation while the answer [No] means
537 that the paper has limitations, but those are not discussed in the paper.
- 538 • The authors are encouraged to create a separate “Limitations” section in their paper.
- 539 • The paper should point out any strong assumptions and how robust the results are to
540 violations of these assumptions (e.g., independence assumptions, noiseless settings,
541 model well-specification, asymptotic approximations only holding locally). The authors
542 should reflect on how these assumptions might be violated in practice and what the
543 implications would be.
- 544 • The authors should reflect on the scope of the claims made, e.g., if the approach was
545 only tested on a few datasets or with a few runs. In general, empirical results often
546 depend on implicit assumptions, which should be articulated.
- 547 • The authors should reflect on the factors that influence the performance of the approach.
548 For example, a facial recognition algorithm may perform poorly when image resolution
549 is low or images are taken in low lighting. Or a speech-to-text system might not be
550 used reliably to provide closed captions for online lectures because it fails to handle
551 technical jargon.
- 552 • The authors should discuss the computational efficiency of the proposed algorithms
553 and how they scale with dataset size.
- 554 • If applicable, the authors should discuss possible limitations of their approach to
555 address problems of privacy and fairness.
- 556 • While the authors might fear that complete honesty about limitations might be used by
557 reviewers as grounds for rejection, a worse outcome might be that reviewers discover
558 limitations that aren’t acknowledged in the paper. The authors should use their best
559 judgment and recognize that individual actions in favor of transparency play an impor-
560 tant role in developing norms that preserve the integrity of the community. Reviewers
561 will be specifically instructed to not penalize honesty concerning limitations.

562 3. Theory assumptions and proofs

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

565
566
567
568
569
570
571
572
573
574
575
576
577
578
579
580
581
582
583
584
585
586
587
588
589
590
591
592
593
594
595
596
597
598
599
600
601
602
603
604
605
606
607
608
609
610
611
612
613
614
615
616
617
618

Answer: [N/A]

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

Guidelines:

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

4. Experimental result reproducibility

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

Answer: [Yes]

Justification: Sections 4 and 5 describe the evaluation protocol and reported metrics, and Appendix A records completed run identifiers, exact dataset revision, code/data availability, and compute details.

Guidelines:

- The answer [N/A] means that the paper does not include experiments.
- If the paper includes experiments, a [No] answer to this question will not be perceived well by the reviewers: Making the paper reproducible is important, regardless of whether the code and data are provided or not.
- If the contribution is a dataset and/or model, the authors should describe the steps taken to make their results reproducible or verifiable.
- Depending on the contribution, reproducibility can be accomplished in various ways. For example, if the contribution is a novel architecture, describing the architecture fully might suffice, or if the contribution is a specific model and empirical evaluation, it may be necessary to either make it possible for others to replicate the model with the same dataset, or provide access to the model. In general, releasing code and data is often one good way to accomplish this, but reproducibility can also be provided via detailed instructions for how to replicate the results, access to a hosted model (e.g., in the case of a large language model), releasing of a model checkpoint, or other means that are appropriate to the research performed.
- While NeurIPS does not require releasing code, the conference does require all submissions to provide some reasonable avenue for reproducibility, which may depend on the nature of the contribution. For example
 - (a) If the contribution is primarily a new algorithm, the paper should make it clear how to reproduce that algorithm.
 - (b) If the contribution is primarily a new model architecture, the paper should describe the architecture clearly and fully.
 - (c) If the contribution is a new model (e.g., a large language model), then there should either be a way to access this model for reproducing the results or a way to reproduce the model (e.g., with an open-source dataset or instructions for how to construct the dataset).
 - (d) We recognize that reproducibility may be tricky in some cases, in which case authors are welcome to describe the particular way they provide for reproducibility. In the case of closed-source models, it may be that access to the model is limited in some way (e.g., to registered users), but it should be possible for other researchers to have some path to reproducing or verifying the results.

619
620
621
622
623
624
625
626
627
628
629
630
631
632
633
634
635
636
637
638
639
640
641
642
643
644
645
646
647
648
649
650
651
652
653
654
655
656
657
658
659
660
661
662
663
664
665
666
667
668
669
670

5. Open access to data and code

Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material?

Answer: [Yes]

Justification: The paper describes the benchmark runtime, scoring code, and task packaging as releasable assets, and the appendix describes the anonymous source artifact, reviewer data mirror, pinned task revision, and result artifacts.

Guidelines:

- The answer [N/A] means that paper does not include experiments requiring code.
- Please see the NeurIPS code and data submission guidelines (<https://neurips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- While we encourage the release of code and data, we understand that this might not be possible, so [No] is an acceptable answer. Papers cannot be rejected simply for not including code, unless this is central to the contribution (e.g., for a new open-source benchmark).
- The instructions should contain the exact command and environment needed to run to reproduce the results. See the NeurIPS code and data submission guidelines (<https://neurips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- The authors should provide instructions on data access and preparation, including how to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- The authors should provide scripts to reproduce all experimental results for the new proposed method and baselines. If only a subset of experiments are reproducible, they should state which ones are omitted from the script and why.
- At submission time, to preserve anonymity, the authors should release anonymized versions (if applicable).
- Providing as much information as possible in supplemental material (appended to the paper) is recommended, but including URLs to data and code is permitted.

6. Experimental setting/details

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

Answer: [Yes]

Justification: Sections 4 and 5 specify the relevant settings for this evaluation paper: provider, model, access mode, reasoning level, task coverage, benchmark score definition, per-difficulty reporting, and build-success aggregation.

Guidelines:

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

7. Experiment statistical significance

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

Answer: [Yes]

Justification: The paper reports bootstrap confidence intervals and exact paired sign-flip tests for repeated full-sweep checks across seven standalone configurations where repeat counts support them. It still explicitly avoids pairwise model-ranking significance claims for configurations without enough repeated independent sweeps, and Appendix A states how repeated-run aggregates should be reported in future leaderboard versions.

Guidelines:

- The answer [N/A] means that the paper does not include experiments.

- 671 • The authors should answer [Yes] if the results are accompanied by error bars, confidence
672 intervals, or statistical significance tests, at least for the experiments that support the
673 main claims of the paper.
- 674 • The factors of variability that the error bars are capturing should be clearly stated (for
675 example, train/test split, initialization, random drawing of some parameter, or overall
676 run with given experimental conditions).
- 677 • The method for calculating the error bars should be explained (closed form formula,
678 call to a library function, bootstrap, etc.)
- 679 • The assumptions made should be given (e.g., Normally distributed errors).
- 680 • It should be clear whether the error bar is the standard deviation or the standard error
681 of the mean.
- 682 • It is OK to report 1-sigma error bars, but one should state it. The authors should
683 preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis
684 of Normality of errors is not verified.
- 685 • For asymmetric distributions, the authors should be careful not to show in tables or
686 figures symmetric error bars that would yield results that are out of range (e.g., negative
687 error rates).
- 688 • If error bars are reported in tables or plots, the authors should explain in the text how
689 they were calculated and reference the corresponding figures or tables in the text.

690 8. Experiments compute resources

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

694 Answer: [Yes]

695 Justification: Appendix A reports host hardware, key software versions, full-sweep wall-
696 clock ranges for standalone and agent rows, and how provider quota failures are accounted
697 for.

698 Guidelines:

- 699 • The answer [N/A] means that the paper does not include experiments.
- 700 • The paper should indicate the type of compute workers CPU or GPU, internal cluster,
701 or cloud provider, including relevant memory and storage.
- 702 • The paper should provide the amount of compute required for each of the individual
703 experimental runs as well as estimate the total compute.
- 704 • The paper should disclose whether the full research project required more compute
705 than the experiments reported in the paper (e.g., preliminary or failed experiments that
706 didn't make it into the paper).

707 9. Code of ethics

708 Question: Does the research conducted in the paper conform, in every respect, with the
709 NeurIPS Code of Ethics [https://neurips.cc/public/EthicsGuidelines?](https://neurips.cc/public/EthicsGuidelines)

710 Answer: [Yes]

711 Justification: The work is a benchmark and evaluation study, does not involve deception or
712 human subjects, and includes explicit discussion of limitations, failure modes, and broader
713 impacts in Sections 6 and 7.

714 Guidelines:

- 715 • The answer [N/A] means that the authors have not reviewed the NeurIPS Code of
716 Ethics.
- 717 • If the authors answer [No], they should explain the special circumstances that require a
718 deviation from the Code of Ethics.
- 719 • The authors should make sure to preserve anonymity (e.g., if there is a special consid-
720 eration due to laws or regulations in their jurisdiction).

721 10. Broader impacts

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

724
725
726
727
728
729
730
731
732
733
734
735
736
737
738
739
740
741
742
743
744
745
746
747
748
749
750
751
752
753
754
755
756
757
758
759
760
761
762
763
764
765
766
767
768
769
770
771
772
773
774
775
776
777

Answer: [Yes]

Justification: The Limitations section and Appendix A discuss positive impacts for engineering productivity and accessibility, and negative impacts including unsafe hardware design, deceptive artifacts, and overtrust in generated designs.

Guidelines:

- The answer [N/A] means that there is no societal impact of the work performed.
- 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 (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.
- The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate Deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
- The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
- If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

11. Safeguards

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

Answer: [N/A]

Justification: The paper releases a benchmark runtime and evaluation assets, not a new high-risk generative model or scraped dataset that would require controlled release safeguards of the type described here.

Guidelines:

- 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, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.
- Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
- 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.

12. Licenses for existing assets

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

Answer: [Yes]

Justification: The bibliography credits the main upstream research assets, and Appendix A includes a dedicated license table for the primary software dependencies, project code, public task dataset, authored fixtures, and provider terms.

778
779
780
781
782
783
784
785
786
787
788
789
790
791
792
793
794
795
796
797
798
799
800
801
802
803
804
805
806
807
808
809
810
811
812
813
814
815
816
817
818
819
820
821
822
823
824
825
826
827
828
829

Guidelines:

- 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 (e.g., CC-BY 4.0) should be included for each asset.
- For scraped data from a particular source (e.g., website), 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 popular datasets, paperswithcode.com/datasets has curated licenses for some datasets. Their licensing guide can help determine the license of a dataset.
- For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.
- If this information is not available online, the authors are encouraged to reach out to the asset’s creators.

13. New assets

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

Answer: [Yes]

Justification: The benchmark task taxonomy, runtime, scoring protocol, run accounting, and release structure are documented in Sections 4–6 and Appendix A.

Guidelines:

- The answer [N/A] means that the paper does not release new assets.
- Researchers should communicate the details of the dataset/code/model as part of their submissions via structured templates. This includes details about training, license, limitations, etc.
- The paper should discuss whether and how consent was obtained from people whose asset is used.
- At submission time, remember to anonymize your assets (if applicable). You can either create an anonymized URL or include an anonymized zip file.

14. Crowdsourcing and research with human subjects

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

Answer: [N/A]

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

Guidelines:

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

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

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

830
831
832
833
834
835
836
837
838
839
840
841
842
843
844
845
846
847
848
849
850
851
852
853
854

Answer: [N/A]

Justification: The paper does not involve human subjects research.

Guidelines:

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

16. Declaration of LLM usage

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

Answer: [N/A]

Justification: LLMs are not involved as any important, original, or non-standard components.

Guidelines:

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