



# From Task-First Evaluation to Skill-Centered Assessment

---

Anonymous Author(s)

Affiliation

Address

email

1 <https://anonyy-coder.github.io/SkillLens/>

## Abstract

2 Agent skills have become a practical way to extend large language model agents,  
3 but the growing skill ecosystem still lacks a reliable way to judge whether a skill is  
4 worth deploying. Existing evaluations remain largely task-first: they assess skills  
5 through performance on predefined tasks in fixed environments. As skill market-  
6 places expand, this paradigm becomes inadequate: fixed suites can blur a skill’s  
7 marginal contribution with backbone strength and miss its value when tasks fall  
8 outside the skill’s intended scope. We introduce **SkillLens**, a framework for *skill-*  
9 *centered assessment*. Rather than starting from a fixed task suite, SkillLens starts  
10 from the skill artifact itself and automatically constructs capability-aligned evalua-  
11 tion tasks from the skill package. It executes these tasks in isolated environments,  
12 collects execution evidence, and combines automated checks with LLM-based  
13 judging to produce auditable measures of utility, efficiency, and risk. This design  
14 turns arbitrary skills into transparent, deployment-relevant reports and better aligns  
15 evaluation with the adoption question a developer actually faces: what does this  
16 skill enable, and what risk does it introduce? To make these signals available at  
17 discovery time, we surface SkillLens results through a browser extension, and to  
18 support reproducibility we will open-source intermediate trajectories and related  
19 evaluation artifacts at <https://github.com/anonyy-coder/SkillLens>.

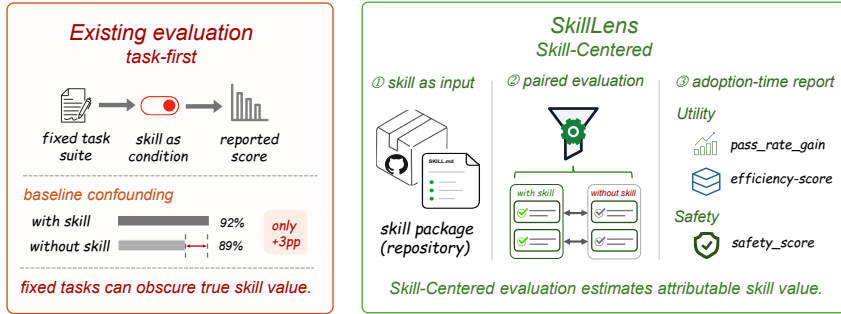
## 20 1 Introduction

21 Agent skills are becoming a common way to extend large language model agents. A skill is a shareable  
22 artifact, usually a directory rooted at SKILL.md and bundled with scripts, templates, dependencies,  
23 and metadata, that an agent can load on demand to specialize to a family of tasks [Anthropic, 2025].  
24 Since Anthropic released the Agent Skills format in late 2025 and OpenAI’s Codex CLI adopted  
25 a structurally similar specification, public ecosystems have expanded from dozens of artifacts to  
26 tens of thousands; a recent survey counts more than 47,000 unique skills across GitHub, community  
27 marketplaces, and corporate contributions [Li et al., 2026]. Skills occupy a distinct layer of the agent  
28 stack. They are neither models nor single tool calls. Instead, they package procedural knowledge,  
29 executable code, and environment assumptions into units that can be installed, audited, and reused.

30 Despite their growing importance, developers still lack a principled way to judge whether a skill  
31 is worth enabling. The practical adoption question is not whether an agent performs well on a  
32 benchmark in the abstract, but whether a particular skill improves task success, reduces effort, or  
33 introduces new safety risk in the workflow it claims to support. Yet GitHub repositories and skill hubs  
34 mostly surface stars, downloads, and README claims rather than direct evidence about a skill’s

**Table 1:** Comparison of SkillLens with representative prior benchmarks across key evaluation properties. Check marks denote supported properties.

Benchmark	Skill artifact	Skill-centered	Utility	Safety	Efficiency	Adoption surfacing
AgentBench Liu et al., 2024	×	×	✓	×	×	×
$\tau$ -bench Yao et al., 2025	×	×	✓	×	×	×
TheAgentCompany Xu et al., 2025	×	×	✓	×	×	×
ToolLLM Qin et al., 2024	×	×	✓	×	×	×
BFCL Patil et al., 2025	×	×	✓	×	×	×
InjecAgent Zhan et al., 2024	×	×	×	✓	×	×
SkillsBench Li et al., 2026	✓	×	✓	×	×	×
<b>SkillLens (ours)</b>	✓	✓	✓	✓	✓	✓



**Figure 1:** SkillLens centers evaluation on each skill artifact rather than benchmarking skills on a fixed task suite, using matched with-skill and no-skill runs to estimate utility and safety.

35 utility and risk. An increasingly important layer of the agent ecosystem therefore remains hard to  
 36 assess at the moment of installation.

37 Existing benchmarks illuminate adjacent issues, but they remain largely *task-first*. AgentBench [Liu  
 38 et al., 2024],  $\tau$ -bench [Yao et al., 2025], and TheAgentCompany [Xu et al., 2025] evaluate agents  
 39 on fixed task suites; ToolLLM [Qin et al., 2024] and BFCL [Patil et al., 2025] study tool invocation  
 40 behavior; InjecAgent [Zhan et al., 2024] emphasizes adversarial robustness; and SkillsBench [Li  
 41 et al., 2026] treats skills as first-class artifacts but still evaluates them on a curated task suite. Table 1  
 42 summarizes this landscape. The shared limitation is that tasks are fixed in advance and a skill is  
 43 judged only through its effect on that suite. For open skill ecosystems, this setup can blur marginal  
 44 skill value with backbone strength and overlook the skill’s intended operating range.

45 We therefore introduce **SkillLens**, a framework for *skill-centered assessment*. Rather than asking how  
 46 a skill moves a fixed benchmark, SkillLens asks how an arbitrary skill package can be systematically  
 47 turned into evaluation tasks that reflect what the skill is actually designed to do. As illustrated in  
 48 Figure 1, the central technical problem is therefore not merely to compare with-skill and no-skill  
 49 runs, but to derive diagnostic tasks from the skill artifact itself.

50 SkillLens addresses this problem through a task-construction pipeline grounded in the skill package.  
 51 It reads `SKILL.md` together with scripts, dependencies, and configuration files, uses an LLM to  
 52 recover the skill’s core capabilities and workflow, and synthesizes three utility schemes that cover the  
 53 skill’s main functions. In parallel, it statically scans the full package to derive security probes. Scheme  
 54 construction enforces a strict separation of roles: the agent sees only business-goal instructions, while  
 55 the evaluator retains the capability analysis, scoring rubric, environment requirements, and expected  
 56 outcomes. To ensure that the resulting scenarios are genuinely discriminative, SkillLens also analyzes  
 57 the behavior boundary of a baseline agent and calibrates tasks around what the baseline can and  
 58 cannot do. It then compiles each scheme into a complete Harbor task artifact, including the required  
 59 inputs, environment setup, and `task.toml`, and instantiates mirrored with-skill and no-skill variants  
 60 for paired evaluation.

61 Once tasks are constructed in this skill-grounded way, the rest of the framework turns execution  
62 into auditable evidence rather than a single opaque score. SkillLens runs each artifact in an isolated  
63 Harbor container, records outputs, full trajectories, filesystem diffs, and outbound network activity,  
64 and uses these traces for automated evaluation. Rule-based checks and LLM judges estimate utility  
65 and efficiency, while static findings over 21 risk patterns trigger targeted dynamic probes that assess  
66 exploitability and produce a safety score. The final report jointly summarizes utility, efficiency, and  
67 safety, with every conclusion traceable to concrete execution records.

68 This design keeps evaluation both skill-grounded and adoption-facing. Because tasks are derived  
69 from a skill’s claimed capabilities rather than from a fixed benchmark, the resulting judgments align  
70 more closely with what developers need to know at installation time, while the underlying trajectories  
71 and process artifacts remain available for inspection.

72 Our contributions are threefold.

- 73 • **A skill-centered evaluation perspective.** We frame skill assessment around the skill artifact  
74 itself, rather than treating a skill merely as a condition within a fixed task suite. This  
75 perspective isolates marginal skill utility and reduces both baseline conflation and task–skill  
76 mismatch.
- 77 • **An evidence-grounded evaluation pipeline for open skill ecosystems.** We introduce  
78 SkillLens, which automatically derives capability-matched utility schemes and security  
79 probes, executes paired tasks in instrumented sandboxes, and turns multimodal execution  
80 evidence into auditable measures of utility, efficiency, and safety.
- 81 • **Controlled and real-world validation with open intermediate artifacts.** We validate  
82 the framework on both controlled safety cases and real-world skill packages across agent  
83 backbones, and we will open-source intermediate trajectories and related process artifacts to  
84 support reproducibility and downstream research.

## 85 2 SkillLens

86 SkillLens evaluates the skill artifact itself. Starting from a collected package, it derives skill-grounded  
87 utility tasks and safety probes, executes them in Harbor, and aggregates the resulting evidence into  
88 deployment-facing reports over utility, efficiency, and safety.

89 We organize the framework into four subsections covering five stages: skill collection, scheme  
90 generation, task construction, sandboxed execution, and automated evaluation and reporting. The  
91 central design goal is to derive tasks from the skill’s claimed capabilities and runtime assumptions  
92 rather than from a fixed benchmark.

### 93 2.1 Skill Collection

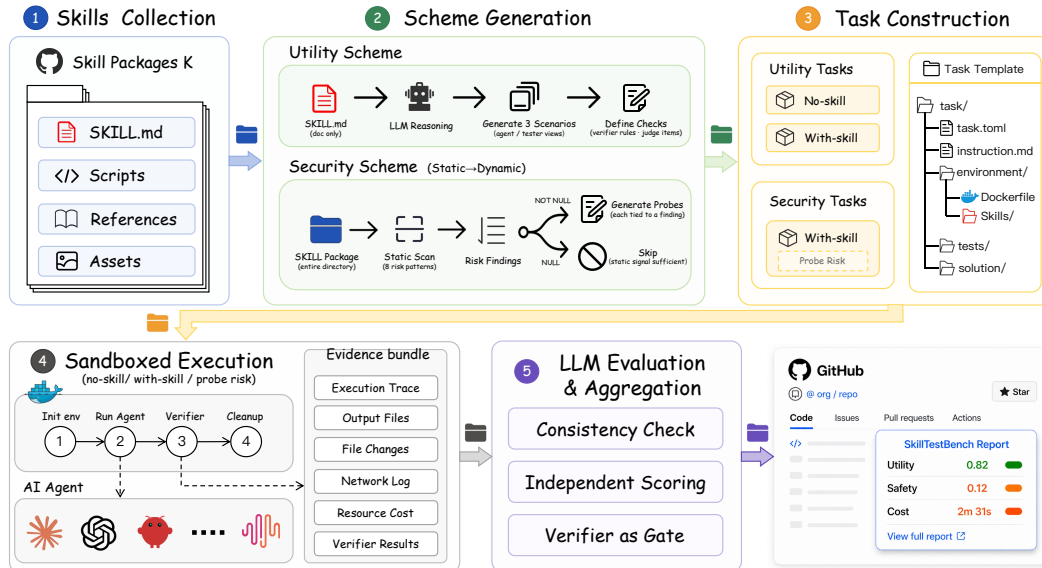
94 The pipeline begins by filtering raw external packages into benchmarkable skills and normalizing each  
95 admitted package into a common three-part representation: the documentation rooted at `SKILL.md`,  
96 the implementation bundle, and provenance metadata such as source, owner, category, and package  
97 name.

98 A package is admitted when it contains a root `SKILL.md`. This matches prevailing skill-package  
99 conventions without assuming a fixed API, programming language, or script layout; those details are  
100 interpreted later during scheme generation, risk scanning, and execution. This stage yields a collected  
101 set of normalized skills, together with the provenance records needed to trace every downstream  
102 scheme, task, execution bundle, and report back to its source package.

### 103 2.2 Scheme Generation

104 Before we build executable tasks in Harbor, our sandboxed evaluation framework, we first derive  
105 two intermediate schemes from each skill. The utility scheme comes from `SKILL.md` and encodes  
106 intended usage scenarios. The security scheme comes from the full package and encodes what risk  
107 paths need to be checked. These schemes determine the tasks built in the next stage.

108 **Utility Scheme.** The utility scheme starts from `SKILL.md` and extracts a small set of representative  
109 usage scenarios. Each scenario records the instruction shown to the agent, the execution context and



**Figure 2: SkillLens pipeline.** A collected skill is converted into utility and security schemes, compiled into utility and security tasks, executed in isolated sandboxes, and aggregated into utility, safety, and resource diagnostics.

110 materials, and the checks used later for evaluation, including verifier rules and judge items. Each  
 111 scenario later becomes a matched no-skill/with-skill task pair.

112 **Security Scheme.** The security scheme starts with a static scan over the full skill package and  
 113 produces a set of risk findings. The scan covers patterns such as prompt injection, credential access,  
 114 unsafe file or network behavior, hidden instructions, supply-chain issues, and robustness failures. If  
 115 the scan finds nothing, the security scheme ends with the static result. If it finds risks, we create  
 116 one probe specification for each finding. Those probe specifications are then carried forward to task  
 117 construction.

### 118 2.3 Task Construction and Sandboxed Execution

119 The generated schemes are then compiled into Harbor tasks. As shown in Figure 2, every task follows  
 120 the same task format: an agent-facing instruction, a sandboxed environment, verifier logic, and  
 121 reference materials. The resulting tasks are then executed under the same instrumented sandboxed-  
 122 execution pipeline, which holds task setup fixed across runs.

123 **Utility Tasks and Security Tasks.** Each utility scenario is compiled into two utility tasks: a no-  
 124 skill task and a with-skill task. The two tasks share the same instruction, inputs, environment, and  
 125 evaluation criteria. The only difference is whether the evaluated skill is exposed through the scenario’s  
 126 contract. This paired construction supports attribution to skill availability.

127 Security tasks are constructed differently. Each probe specification becomes one with-skill risk task  
 128 tied to a single finding. Here the goal is not counterfactual comparison, but runtime confirmation of a  
 129 specific risk path when the skill is present. Utility tasks and security tasks still reuse the same task  
 130 template and sandboxed execution pipeline, so they share tooling, logs, and provenance.

131 **Shared Evidence.** Each trial yields a common evidence bundle containing the execution trace, output  
 132 files, file changes, network log, resource cost, and verifier results—the same evidence bundle shown  
 133 in Figure 2. A utility scenario therefore yields two evidence bundles, one from the no-skill run and  
 134 one from the with-skill run, while a security task yields one risk-oriented bundle. Using one evidence  
 135 format for semantic judgment, safety confirmation, and cost accounting keeps all reported numbers  
 136 traceable to concrete executions.

137 **2.4 Aggregation and Reporting**

138 Aggregation turns the shared execution records into three reported axes: utility, safety, and resource  
 139 cost. Utility comes from matched no-skill/with-skill comparisons, safety comes from static findings  
 140 and runtime confirmation, and cost comes from the same execution records. The public report surfaces  
 141 utility and safety directly, while cost is reported through explicit resource fields in the diagnostics.

142 **Utility Scoring.** For each skill, let  $\mathcal{V}_j$  denote the set of matched pairs that support attribution: both  
 143 runs admit semantic judgment, and the no-skill run does not receive access to the evaluated skill.  
 144 Pairs that fail these checks are retained in the diagnostics rather than silently dropped. For each valid  
 145 pair, the LLM judge scores the no-skill and with-skill outputs against the same semantic checklist.  
 146 Let  $n_{j,k}^{\text{wi}}$  and  $n_{j,k}^{\text{wo}}$  be the numbers of passed judge items out of  $N_{j,k}$  total items. We report utility as  
 147 the mean pass-rate gain across valid pairs:

$$U_{\text{pass}}(\mathcal{K}_j) = \frac{1}{|\mathcal{V}_j|} \sum_{(j,k) \in \mathcal{V}_j} \frac{\max(n_{j,k}^{\text{wi}} - n_{j,k}^{\text{wo}}, 0)}{N_{j,k}}, \quad (1)$$

148 Negative differences are clipped at zero because  $U_{\text{pass}}(\mathcal{K}_j)$  is intended to measure realized marginal  
 149 gain from enabling the skill, while regressions remain visible in the per-scenario diagnostics.

150 We report efficiency separately on the subset  $\mathcal{V}_j^{\text{eff}} \subseteq \mathcal{V}_j$  for which both runs pass at least one judge  
 151 item. Let  $\tilde{q} = q_{\text{input}} - q_{\text{cache}}$  denote effective input tokens. For each such pair, we compute relative  
 152 savings in wall-clock time and effective input tokens,

$$e_{j,k}^{(t)} = \frac{t_{j,k}^{\text{wo}} - t_{j,k}^{\text{wi}}}{t_{j,k}^{\text{wo}}}, \quad e_{j,k}^{(q)} = \frac{\tilde{q}_{j,k}^{\text{wo}} - \tilde{q}_{j,k}^{\text{wi}}}{\tilde{q}_{j,k}^{\text{wo}}}, \quad (2)$$

153 and define the per-pair efficiency score as  $\text{efficiency\_score} = (e_{j,k}^{(t)} + e_{j,k}^{(q)})/2$ . We then report

$$U_{\text{eff}}(\mathcal{K}_j) = \frac{1}{|\mathcal{V}_j^{\text{eff}}|} \sum_{(j,k) \in \mathcal{V}_j^{\text{eff}}} \frac{e_{j,k}^{(t)} + e_{j,k}^{(q)}}{2}. \quad (3)$$

154 If no pair is valid, both utility fields are reported as null; if no pair supports efficiency scoring,  
 155 only  $U_{\text{eff}}(\mathcal{K}_j)$  is null. Verifier outcomes only determine whether a pair is structurally comparable;  
 156 semantic success still comes from the judge items.

157 **Resource Cost.** For interpretability, we also retain the raw resource deltas over the same efficiency-  
 158 valid pairs:

$$\Delta t(\mathcal{K}_j) = \frac{1}{|\mathcal{V}_j^{\text{eff}}|} \sum_{(j,k) \in \mathcal{V}_j^{\text{eff}}} (t_{j,k}^{\text{wi}} - t_{j,k}^{\text{wo}}), \quad \Delta \tilde{q}(\mathcal{K}_j) = \frac{1}{|\mathcal{V}_j^{\text{eff}}|} \sum_{(j,k) \in \mathcal{V}_j^{\text{eff}}} (\tilde{q}_{j,k}^{\text{wi}} - \tilde{q}_{j,k}^{\text{wo}}). \quad (4)$$

159 When  $\mathcal{V}_j^{\text{eff}} = \emptyset$ , these raw resource fields are also reported as null. These cost fields are not folded  
 160 into a single public score; instead, they are surfaced as explicit diagnostics alongside utility and  
 161 safety.

162 **Safety Scoring.** The scanner from Section 2.2 produces a set of findings. Each finding has a severity  
 163 level, an existence confidence, and, when a dynamic probe is available, an exploitability estimate.  
 164 Severity is mapped to a fixed base penalty with  $w(H) = 15$ ,  $w(M) = 10$ , and  $w(L) = 5$ .

165 When a dynamic probe is available, a security judge estimates an exploitability confidence  $c_{j,\ell}^{\text{exploit}}$   
 166 from verifier signals, trajectory evidence, and the probe’s expected runtime signature; otherwise we  
 167 use a default value of 0.6 and mark that case in the diagnostics. Appendix Table A4 summarizes  
 168 how dynamic verdicts map to exploitability ranges. The `agent_refused` case is not mapped to zero,  
 169 because the risky path may still exist in the package even if one agent refuses to execute it.

170 The final safety score starts from 100 and subtracts the confidence-weighted penalties from all  
 171 findings:

$$S(\mathcal{K}_j) = \max\left(10, 100 - \sum_{\ell=1}^m w(\sigma_{j,\ell}) c_{j,\ell}^{\text{exist}} c_{j,\ell}^{\text{exploit}}\right). \quad (5)$$

172 The lower bound preserves score separation among risky skills while keeping the scale interpretable.  
 173 We report Pass when  $S(\mathcal{K}_j) = 100$ , Caution when  $80 \leq S(\mathcal{K}_j) < 100$ , and Risky otherwise,  
 174 following the threshold convention in SkillTester [Wang et al., 2026]. This score is intended as a  
 175 deployment-facing risk summary rather than as a replacement for the underlying finding records,  
 176 which remain available in the diagnostics.

177 **Final Report.** The final report presents the public fields side by side:

$$\text{Report}(\mathcal{K}_j) = (U_{\text{pass}}(\mathcal{K}_j), U_{\text{eff}}(\mathcal{K}_j), S(\mathcal{K}_j), \text{diag}(\mathcal{K}_j)), \quad (6)$$

178 where the public summary surfaces attributable effectiveness, normalized efficiency, and safety, while  
 179  $\text{diag}(\mathcal{K}_j)$  retains the supporting details: invalid pairs and access-pattern flags, raw resource deltas  
 180 such as  $\Delta t$  and  $\Delta \bar{q}$ , per-scenario time/token sub-scores, per-scenario regressions, and per-finding  
 181 safety records with default markers. Presenting these fields side by side makes disagreement across  
 182 utility, safety, and resource cost visible instead of hiding it behind a single score.

183 **Browser Extension.** To expose these results at the point of adoption, we build a lightweight browser  
 184 extension for skill-hosting platforms. It links candidate repositories and marketplace entries to their  
 185 SkillLens report pages, surfacing pass-rate gain, efficiency, safety score, and supporting diagnostics  
 186 in a unified view. The extension therefore turns SkillLens from an offline evaluation pipeline into a  
 187 discovery-time decision aid.

### 188 3 Empirical Evaluation of SkillLens

189 This section evaluates SkillLens on **226 real-world skill packages** collected from public skill  
 190 repositories and community-maintained skill indexes, all normalized into the representation defined in  
 191 Section 2.1. For distributional analyses, generated utility scenarios are grouped into **23 occupational**  
 192 **categories**. Across the experiments, we use matched no-skill and with-skill executions so that utility  
 193 and safety are attributed under comparable conditions: utility is evaluated across six agent–model  
 194 configurations via PRG and  $U_{\text{eff}}$ , while dynamic safety is evaluated on the three representative  
 195 configurations that support probe instrumentation— Claude Code / Sonnet 4.6, Codex / GPT-5.1,  
 196 and Codex / GPT-5.4—via  $S$ . Detailed scanner statistics, dynamic-probe diagnostics, and full  
 197 configuration settings are deferred to the appendices.

#### 198 3.1 Utility across Agent–Model Configurations

199 Unless noted otherwise, utility attribution is computed under Policy A independently within each  
 200 configuration. For each utility scenario, we execute matched no-skill and with-skill conditions once,  
 201 holding the configuration, task instance, and, when supported, random seed fixed. We exclude cases  
 202 in which the with-skill run fails because of a timeout or judge failure, but retain cases in which the  
 203 no-skill baseline times out and the with-skill run completes; after filtering, each configuration retains  
 204 between **584 and 676** scenarios, and the Codex / GPT-5.4 configuration retains 653.

205 Table 2 provides a single-configuration view under Codex / GPT-5.4 and makes SkillLens’s three  
 206 deployment-facing axes concrete. Mean with-skill pass rate is 0.944, but PRG still spans 0.060–  
 207 0.407, from *building-grounds-cleaning* to *educational-instruction*. Mean efficiency is moderate  
 208 ( $U_{\text{eff}} = 0.211$ ), with 9 categories running faster with the skill than without it, most clearly in  
 209 *transportation-material-moving* (195 s vs. 243 s). Safety is likewise uneven: the overall score is  
 210 95.0, but *architecture-engineering* and *personal-care-service* together account for 39.8% of all  
 211 high-severity static findings, while the Pearson correlation between PRG and security score is only  
 212 0.111.

213 Figure 3 reports category-level mean PRG across the six configurations. The top row gives each  
 214 configuration’s overall mean PRG, and the column headers report the corresponding no-skill baseline  
 215 pass rate (wo). Three findings are most important. **Cross-configuration mean PRG decreases**  
 216 **as the no-skill baseline strengthens.** Across the six configurations, mean PRG and mean wo are  
 217 strongly negatively correlated ( $r = -0.93$ ). Configurations with weaker baselines (Codex / GPT-5.1,  
 218 wo = 0.608; Claude Code / Sonnet 4, wo = 0.636) report the highest mean PRG (0.288 and 0.272),  
 219 whereas configurations with stronger baselines (Claude Code / Sonnet 4.6, wo = 0.776; OpenCode  
 220 / Sonnet 4.6, wo = 0.771) report the lowest (0.185 and 0.171). This pattern is consistent with a  
 221 ceiling effect: stronger no-skill baselines leave less residual headroom for a skill to unlock additional

**Table 2:** Utility and security metrics by occupation category (codex + gpt-5.4, Policy A).  $w_i$  = with skill;  $w_o$  = without skill. PRG = pass rate gain =  $\max(0, (w_i\_passed - w_o\_passed)/total\_items)$ . Eff. tokens = effective input tokens (input – cache tokens), averaged per scenario. Time/token/efficiency scores  $\in [0, 1]$ . Static findings H/M/L = **total count** across all skills in the category (not per-skill average). Security score  $\in [0, 100]$ .

Category	Pass Rate			Time (s)		Eff. Tokens		Efficiency Score			Sec.	Static Findings		
	$w_i$	$w_o$	PRG	$w_i$	$w_o$	$w_i$	$w_o$	Time	Token	Eff.	Score	H	M	L
educational-instruction	0.911	0.516	0.407	164	121	54K	35K	0.089	0.226	0.158	90.6	7	3	5
community-social-service	0.944	0.581	0.369	118	101	12K	10K	0.076	0.161	0.118	96.3	3	2	5
office-administrative	0.936	0.632	0.309	130	175	49K	61K	0.337	0.497	0.417	92.9	7	0	4
management	0.969	0.680	0.289	184	168	38K	43K	0.161	0.247	0.204	99.7	0	1	1
computer-mathematical	0.902	0.648	0.271	108	113	36K	34K	0.225	0.419	0.322	96.6	3	3	1
arts-design-media	0.976	0.736	0.241	109	94	15K	14K	0.064	0.267	0.165	99.6	0	0	2
military-specific	0.935	0.737	0.237	163	152	23K	26K	0.176	0.244	0.210	97.0	1	5	1
transportation-material-moving	0.936	0.719	0.228	194	242	43K	48K	0.260	0.284	0.272	96.4	1	2	10
sales-related	0.963	0.754	0.214	171	140	41K	21K	0.033	0.250	0.141	99.3	1	1	2
protective-service	0.912	0.713	0.210	140	120	45K	31K	0.116	0.218	0.167	95.5	2	4	3
food-preparation-serving	0.885	0.690	0.200	116	90	20K	15K	0.108	0.124	0.116	95.6	2	5	8
healthcare-practitioners	0.955	0.763	0.193	201	203	68K	73K	0.128	0.325	0.226	99.8	0	0	1
personal-care-service	0.911	0.756	0.183	136	151	23K	27K	0.186	0.384	0.285	82.2	18	5	9
life-physical-social-science	0.942	0.787	0.182	292	323	104K	92K	0.212	0.167	0.189	91.3	4	8	7
business-financial	0.982	0.800	0.182	212	192	42K	39K	0.062	0.182	0.122	99.1	0	1	3
legal	0.984	0.813	0.176	266	259	120K	123K	0.153	0.315	0.234	98.1	0	6	4
construction-extraction	0.914	0.779	0.141	179	191	41K	33K	0.235	0.250	0.243	98.9	0	2	3
architecture-engineering	0.982	0.868	0.120	167	162	33K	53K	0.168	0.335	0.251	78.9	15	2	15
installation-maint-repair	0.933	0.815	0.118	163	146	47K	31K	0.084	0.292	0.188	99.3	0	1	1
farming-fishing-forestry	0.973	0.871	0.107	167	207	54K	68K	0.273	0.395	0.334	98.1	0	3	1
production	0.937	0.864	0.090	167	152	46K	32K	0.101	0.197	0.149	86.5	10	1	7
healthcare-support	0.954	0.880	0.082	155	117	39K	35K	0.121	0.157	0.139	95.2	5	3	2
building-grounds-cleaning	0.963	0.913	0.060	119	123	21K	18K	0.166	0.313	0.239	96.0	4	1	4
<b>Overall</b>	<b>0.944</b>	<b>0.753</b>	<b>0.200</b>	<b>166</b>	<b>162</b>	<b>44K</b>	<b>42K</b>	<b>0.152</b>	<b>0.270</b>	<b>0.211</b>	<b>95.0</b>	<b>83</b>	<b>59</b>	<b>99</b>

222 passes. Cross-configuration differences in PRG should therefore be interpreted relative to  $w_o$ , not as  
 223 a standalone measure of skill quality.

224 **The same ceiling effect appears within shared harnesses.** The same ordering holds when the  
 225 harness is fixed. Codex / GPT-5.1 (0.288) exceeds Codex / GPT-5.4 (0.200), Claude Code / Sonnet  
 226 4 (0.272) exceeds Claude Code / Sonnet 4.6 (0.185), and OpenCode / GPT-5.4 (0.213) exceeds  
 227 OpenCode / Sonnet 4.6 (0.171). This within-harness ordering again suggests that stronger backbones  
 228 compress observable skill-induced pass-rate gains by reducing available headroom.

229 **Categories and harnesses contribute additional structured variation.** Some categories remain  
 230 consistently high across configurations: *educational-instruction*, *community-social-service*, and  
 231 *computer-mathematical* all sustain high PRG across all six configurations, whereas *building-grounds-*  
 232 *cleaning* and *farming-fishing-forestry* remain low across all six. Mid-range categories, such as  
 233 *office-administrative* and *transportation-material-moving*, reorder substantially across configurations.  
 234 The harness itself also contributes independent variation: under the same backbone, PRG on *office-*  
 235 *administrative* differs by 0.25 between Claude Code / Sonnet 4.6 (0.452) and OpenCode / Sonnet  
 236 4.6 (0.200), indicating that skill-harness compatibility is task-type dependent. Taken together, these  
 237 patterns indicate that PRG should be interpreted relative to baseline headroom and execution context,  
 238 rather than in isolation.

239 Appendix C.1–C.3 provide the ceiling-effect decomposition, Codex / GPT-5.4 baseline-stratified  
 240 results, and within-category scenario-level PRG distributions.

### 241 3.2 Static Scanner Validation

242 We validate the static security scanner introduced in Section 2.4 against expert-labeled findings on  
 243 a separate controlled subset of 29 skills comprising 186 ground-truth findings, disjoint from the  
 244 226-skill end-to-end evaluation. This controlled subset serves as a measurement-validation step:  
 245 before using the scanner for the large-scale safety profile on the 226-skill set in Section 3.3, we first  
 246 test whether it can recover expert-injected risks. This helps ensure that the subsequent ecosystem-level  
 247 safety profile reflects the skill ecosystem rather than artifacts of the detection layer.

248 Table 3 shows that the scanner is intentionally recall-oriented on controlled human-injected risks,  
 249 recovering 90.9% of all findings and maintaining especially high recall on High- and Medium-severity

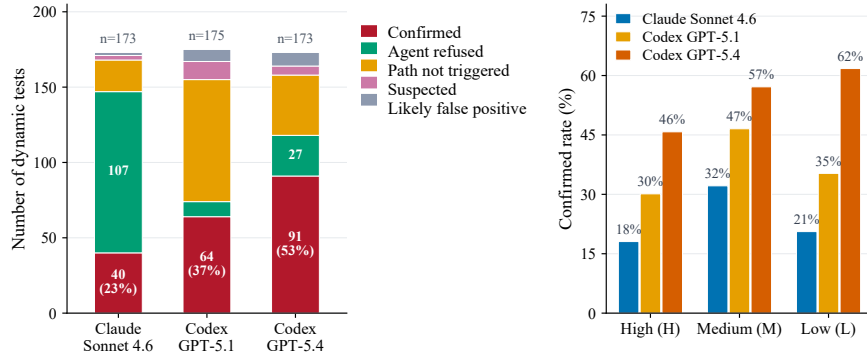


**Figure 3: Mean PRG across configurations and occupations.** Cells report category-level mean PRG over 23 occupational categories and six agent–model configurations. The top row gives the configuration-level mean, and column headers give the no-skill baseline pass rate (wo). Policy A is applied independently per configuration ( $n = 584\text{--}676$ ).

**Table 3: Static-scanner validation on controlled findings.** Detection performance on 29 skills with 186 expert-labeled, human-injected findings. Recall is the primary signal; precision and FI are reported for completeness. #G and #P denote ground-truth and scanner-reported counts.

	#G	#P	TP	FP	FN	P (%)	R (%)	FI (%)
<b>Overall</b>	<b>186</b>	<b>219</b>	<b>169</b>	<b>50</b>	<b>17</b>	<b>77.2</b>	<b>90.9</b>	<b>83.5</b>
High	87	104	84	20	3	80.8	96.5	88.0
Medium	47	60	46	14	1	76.7	97.9	86.0
Low	52	55	39	16	13	70.9	75.0	72.9

250 cases (96.5% and 97.9%). Its main limitation is the expected recall–precision trade-off: the scanner  
 251 over-reports some findings, especially among lower-severity patterns. For a deployment-oriented  
 252 gate, however, missing severe risks is costlier than surfacing some spurious warnings. The detailed  
 253 per-pattern breakdown is reported in Appendix F.



**Figure 4: Dynamic safety outcomes across agent backbones.** Dynamic probes are run on the same static-finding set under three representative agent backbones. Both the verdict distribution and the confirmed rate shift substantially across agents, showing that runtime exploitability is agent-conditioned.

### 254 3.3 Safety Profile of Evaluated Skills

255 We apply the scanner validated in Section 3.2, together with the dynamic probes of Section 2.4, to  
 256 the same 226-skill set evaluated in this section. Dynamic probes are run on the three representative  
 257 configurations that support probe instrumentation: Claude Code / Sonnet 4.6, Codex / GPT-5.1,  
 258 and Codex / GPT-5.4. This produces two complementary views of safety: a package-level score  
 259 distribution and the agent-conditioned dynamic outcomes in Figure 4. The former shows how risk  
 260 is distributed across the ecosystem, whereas the latter reveals how a fixed set of static findings  
 261 materializes into different runtime verdicts under representative agent backbones.

262 **Ecosystem-level distribution.** At the ecosystem level, safety risk is highly concentrated rather than  
 263 uniform. Of the 226 evaluated skills, 157 receive the maximum score, while the remaining 69 form  
 264 a long risky tail that spans a broad severity range, including 18 skills with scores below 70. This skew  
 265 makes a per-skill status gate at  $\theta_s = 80$  more informative than reporting a single corpus average:  
 266 an average would be dominated by the safe majority and would mask the smaller but operationally  
 267 important subset of risky packages.

268 **Agent-conditioned exploitability.** Figure 4 further shows that runtime risk depends materially on  
 269 the serving agent, even when the static finding set is held fixed. The confirmed rate increases from  
 270 roughly 23% under Claude Code / Sonnet 4.6, to 37% under Codex / GPT-5.1, and to 53% under  
 271 Codex / GPT-5.4, while the mass assigned to `agent_refused` shifts in the opposite direction. This  
 272 pattern suggests that refusal should not be interpreted as evidence that no risky path exists; rather, it  
 273 reflects the interaction between the package’s latent exposure and the current agent’s execution policy.  
 274 These results motivate keeping `existence_confidence` and `exploitability` as separate terms  
 275 in Eq. 5: static analysis estimates whether a risky path is present, whereas dynamic probes estimate  
 276 whether that path is actually exercised under a specific agent configuration. Appendix E extends this  
 277 analysis with risk-pattern-, occupation-, and agent-resolved decompositions.

## 278 4 Conclusion

279 We introduced **SkillLens**, a skill-centered framework for evaluating agent skill packages through  
 280 matched with-skill and no-skill runs. Across 226 real-world skills, utility varies with baseline  
 281 headroom, occupational category, and agent-model configuration, while safety remains distinct  
 282 from utility: the scanner attains high recall on controlled findings and dynamic exploitability shifts  
 283 across backbones. These results motivate reporting pass-rate gain, efficiency, and safety as separate  
 284 adoption-time axes. We surface this view through a browser extension; future work will broaden risk  
 285 coverage and calibration.

## 286 References

287 Anthropic. Equipping agents for the real world with Agent Skills. Anthropic En-  
 288 gineering Blog, October 2025. URL <https://www.anthropic.com/engineering/>

- 289 equipping-agents-for-the-real-world-with-agent-skills. Published October 16,  
290 2025; Agent Skills open standard update published December 18, 2025.
- 291 Carlos E. Jimenez, John Yang, Alexander Wettig, Shunyu Yao, Kexin Pei, Ofir Press, and Karthik R.  
292 Narasimhan. SWE-bench: Can language models resolve real-world github issues? In *The Twelfth*  
293 *International Conference on Learning Representations*, 2024. URL [https://openreview.net/](https://openreview.net/forum?id=VTF8yNQm66)  
294 [forum?id=VTF8yNQm66](https://openreview.net/forum?id=VTF8yNQm66). Oral presentation.
- 295 Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman  
296 Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, Sebastian Riedel,  
297 and Douwe Kiela. Retrieval-augmented generation for knowledge-intensive NLP tasks. In  
298 *Advances in Neural Information Processing Systems*, volume 33, pages 9459–9474. Cur-  
299 ran Associates, Inc., 2020. URL [https://proceedings.neurips.cc/paper/2020/hash/](https://proceedings.neurips.cc/paper/2020/hash/6b493230205f780e1bc26945df7481e5-Abstract.html)  
300 [6b493230205f780e1bc26945df7481e5-Abstract.html](https://proceedings.neurips.cc/paper/2020/hash/6b493230205f780e1bc26945df7481e5-Abstract.html).
- 301 Xiangyi Li, Wenbo Chen, Yimin Liu, Shenghan Zheng, Xiaokun Chen, Yifeng He, Yubo Li, Bingran  
302 You, Haotian Shen, Jiankai Sun, Shuyi Wang, Binxu Li, Qunhong Zeng, Di Wang, Xuandong Zhao,  
303 Yuanli Wang, Roey Ben Chaim, Zonglin Di, Yipeng Gao, Junwei He, Yizhuo He, Liqiang Jing,  
304 Luyang Kong, Xin Lan, Jiachen Li, Songlin Li, Yijiang Li, Yueqian Lin, Xinyi Liu, Xuanqing Liu,  
305 Haoran Lyu, Ze Ma, Bowei Wang, Runhui Wang, Tianyu Wang, Wengao Ye, Yue Zhang, Hanwen  
306 Xing, Yiqi Xue, Steven Dillmann, and Han-chung Lee. SkillsBench: Benchmarking how well  
307 agent skills work across diverse tasks, 2026. URL <https://arxiv.org/abs/2602.12670>.
- 308 Xiao Liu, Hao Yu, Hanchen Zhang, Yifan Xu, Xuanyu Lei, Hanyu Lai, Yu Gu, Hangliang Ding,  
309 Kaiwen Men, Kejuan Yang, Shudan Zhang, Xiang Deng, Aohan Zeng, Zhengxiao Du, Chenhui  
310 Zhang, Sheng Shen, Tianjun Zhang, Yu Su, Huan Sun, Minlie Huang, Yuxiao Dong, and Jie Tang.  
311 AgentBench: Evaluating LLMs as agents. In *The Twelfth International Conference on Learning*  
312 *Representations*, 2024. URL <https://openreview.net/forum?id=zAdUB0aCTQ>.
- 313 Shishir G. Patil, Huanzhi Mao, Fanjia Yan, Charlie Cheng-Jie Ji, Vishnu Suresh, Ion Stoica, and  
314 Joseph E. Gonzalez. The berkeley function calling leaderboard (BFCL): From tool use to agentic  
315 evaluation of large language models. In *Proceedings of the 42nd International Conference on*  
316 *Machine Learning*, volume 267 of *Proceedings of Machine Learning Research*, pages 48371–48392.  
317 PMLR, 2025. URL <https://proceedings.mlr.press/v267/patil125a.html>.
- 318 Yujia Qin, Shihao Liang, Yining Ye, Kunlun Zhu, Lan Yan, Yaxi Lu, Yankai Lin, Xin Cong, Xiangru  
319 Tang, Bill Qian, Sihan Zhao, Lauren Hong, Runchu Tian, Ruobing Xie, Jie Zhou, Mark Gerstein,  
320 Dahai Li, Zhiyuan Liu, and Maosong Sun. ToolLLM: Facilitating large language models to master  
321 16000+ real-world APIs. In *The Twelfth International Conference on Learning Representations*,  
322 2024. URL <https://openreview.net/forum?id=dHng200Jjr>. Spotlight.
- 323 Jiawen Shi, Zenghui Yuan, Guiyao Tie, Pan Zhou, Neil Zhenqiang Gong, and Lichao Sun. Prompt  
324 injection attack to tool selection in LLM agents. In *Network and Distributed System Security Sym-*  
325 *posium (NDSS)*, 2026. doi: 10.14722/ndss.2026.230675. URL [https://www.ndss-symposium.](https://www.ndss-symposium.org/ndss-paper/prompt-injection-attack-to-tool-selection-in-llm-agents/)  
326 [org/ndss-paper/prompt-injection-attack-to-tool-selection-in-llm-agents/](https://www.ndss-symposium.org/ndss-paper/prompt-injection-attack-to-tool-selection-in-llm-agents/).
- 327 Theodore R. Sumers, Shunyu Yao, Karthik Narasimhan, and Thomas L. Griffiths. Cognitive ar-  
328 chitectures for language agents. *Transactions on Machine Learning Research*, 2024. URL  
329 <https://openreview.net/forum?id=1i6ZCvf1QJ>.
- 330 Leye Wang, Zixing Wang, and Anjie Xu. SkillTester: Benchmarking utility and security of agent  
331 skills, 2026. URL <https://arxiv.org/abs/2603.28815>.
- 332 Frank F. Xu, Yufan Song, Boxuan Li, Yuxuan Tang, Kritanjali Jain, Mengxue Bao, Zora Zhiruo Wang,  
333 Xuhui Zhou, Zhitong Guo, Murong Cao, Mingyang Yang, Hao Yang Lu, Amaad Martin, Zhe Su,  
334 Leander Melroy Maben, Raj Mehta, Wayne Chi, Lawrence Keunho Jang, Yiqing Xie, Shuyan Zhou,  
335 and Graham Neubig. TheAgentCompany: Benchmarking LLM agents on consequential real world  
336 tasks. In *The Thirty-ninth Annual Conference on Neural Information Processing Systems Datasets*  
337 *and Benchmarks Track*, 2025. URL <https://openreview.net/forum?id=LZnKNApvhG>.
- 338 Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik R. Narasimhan, and Yuan  
339 Cao. ReAct: Synergizing reasoning and acting in language models. In *The Eleventh International*  
340 *Conference on Learning Representations*, 2023. URL [https://openreview.net/forum?id=](https://openreview.net/forum?id=WE_vluYUL-X)  
341 [WE\\_vluYUL-X](https://openreview.net/forum?id=WE_vluYUL-X).

- 342 Shunyu Yao, Noah Shinn, Pedram Razavi, and Karthik R. Narasimhan.  $\tau$ -bench: A benchmark for  
343 tool-agent-user interaction in real-world domains. In *The Thirteenth International Conference on*  
344 *Learning Representations*, 2025. URL <https://openreview.net/forum?id=roNSXZpUDN>.
- 345 Qiusi Zhan, Zhixiang Liang, Zifan Ying, and Daniel Kang. InjecAgent: Benchmarking indirect  
346 prompt injections in tool-integrated large language model agents. In *Findings of the Asso-*  
347 *ciation for Computational Linguistics: ACL 2024*, pages 10471–10506, Bangkok, Thailand,  
348 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.findings-acl.624. URL  
349 <https://aclanthology.org/2024.findings-acl.624/>.

## 350 A Related Work

351 Table 1 provides a compact property-level comparison with representative prior benchmarks. This  
352 appendix expands that comparison along four lines of related work: benchmarking LLM agents,  
353 evaluating tool use and procedural augmentation, benchmarking skills as artifacts, and studying the  
354 security of agent ecosystems.

355 **Benchmarking LLM agents.** A large body of work evaluates agents on fixed task suites. Agent-  
356 Bench [Liu et al., 2024],  $\tau$ -bench [Yao et al., 2025], TheAgentCompany [Xu et al., 2025], and  
357 SWE-bench [Jimenez et al., 2024] differ in domain and emphasis, but they share the same unit of  
358 evaluation: an *agent* tested on pre-specified tasks. Our concern is adjacent but different: not how  
359 strong an agent is on a fixed benchmark, but how much a particular *skill package* changes what that  
360 agent can reliably do.

361 **Tool use and procedural augmentation.** Another line of work studies how models use external  
362 capabilities. ToolLLM/ToolBench [Qin et al., 2024] and BFCL [Patil et al., 2025] evaluate whether  
363 models can correctly invoke tools, while work on retrieval-augmented generation [Lewis et al., 2020],  
364 ReAct-style reasoning [Yao et al., 2023], and agent architectures such as CoALA [Sumers et al.,  
365 2024] explores ways to augment agent behavior. These studies establish that augmentation matters,  
366 but they are mostly concerned with mechanisms of tool use or reasoning rather than with evaluating a  
367 concrete skill package as a deployable artifact.

368 **Skills as first-class evaluation artifacts.** The closest prior work is SkillsBench [Li et al., 2026],  
369 which, to our knowledge, is the first benchmark to explicitly treat skills as evaluation artifacts. We  
370 share that starting point, but differ in evaluation framing. SkillsBench studies curated skills on a  
371 curated task suite, whereas our setting is the open skill ecosystem, where the central question is  
372 how arbitrary real-world skills should be assessed at adoption time. The two works are therefore  
373 complementary rather than competitive.

374 **Security evaluation of agent ecosystems.** A parallel literature studies the security of tool-  
375 augmented agents. InjecAgent [Zhan et al., 2024] and later attack studies such as ToolHijacker [Shi  
376 et al., 2026] focus on prompt injection, tool poisoning, and related robustness failures. This line of  
377 work is essential, but it is centered on adversarial failure modes rather than the routine deployment  
378 question of whether a seemingly benign skill should be enabled. Capability benchmarks, by contrast,  
379 often leave security outside the evaluation scope. Our work is motivated by the need to consider both  
380 when assessing skills for real-world use.

## 381 B Reproducibility—Experimental Configuration Details

382 This appendix reports the experimental configuration details used in Section 3. Items already  
383 summarized at the start of Section 3, including the skill set size, sources, occupational partition, and  
384 reported metrics, are not repeated here.

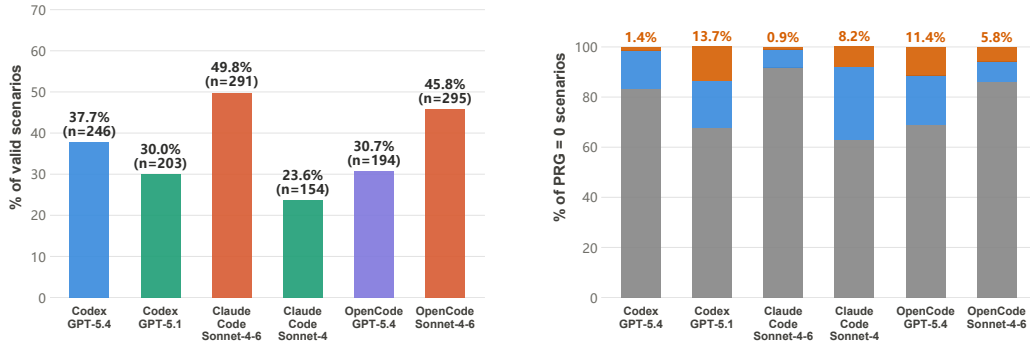
**Table A1:** Detailed experimental configuration for the end-to-end evaluation in Section 3.

Component	Setting
Agent-model configurations	Codex / GPT-5.4, Codex / GPT-5.1, Claude Code / Sonnet 4.6, Claude Code / Sonnet 4, OpenCode / GPT-5.4, OpenCode / Sonnet 4.6
Utility scenarios per skill	3
Valid utility scenarios analyzed	584–676 per configuration under Policy A; 653 for Codex / GPT-5.4 (single-configuration analysis)
Trials per condition (wi / wo)	1
Judge model	Claude Sonnet 4.6 (shared across all configurations)
Controlled scanner-validation set	29 skills, 186 expert-labeled findings

385 **C Supplementary Utility Analyses**

386 **C.1 Ceiling Effect across Agent–Model Configurations**

387 We report the ceiling-effect decomposition used to interpret cross-configuration PRG differences in  
 388 Figure 3. A *strict ceiling* scenario is one in which the no-skill baseline already passes every judge  
 389 item ( $wo = 1.0$ ), forcing  $PRG = 0$  regardless of the skill. We further partition zero-PRG scenarios  
 390 into strict ceiling, *partial ceiling* ( $0.5 \leq wo < 1.0$ ), and *low-baseline no-gain* ( $wo < 0.5$ ). Figure A1  
 391 and Table A2 show that strict ceiling dominates the zero-PRG mass across all configurations. The  
 392 zero-PRG mass therefore carries different meanings across configurations: on stronger baselines it  
 393 predominantly reflects mechanical saturation, whereas on weaker baselines it more often reflects skill  
 394 applicability or skill–agent compatibility.



**Figure A1: Ceiling effect across agent–model configurations.** (a) Share of valid scenarios with  $wo = 1.0$ . (b) Decomposition of zero-PRG scenarios into strict ceiling, partial ceiling, and low-baseline no-gain cases.

**Table A2: Ceiling-effect decomposition under Policy A.** Breakdown columns are percentages within the zero-PRG subset. *LB no-gain* abbreviates low-baseline no-gain.

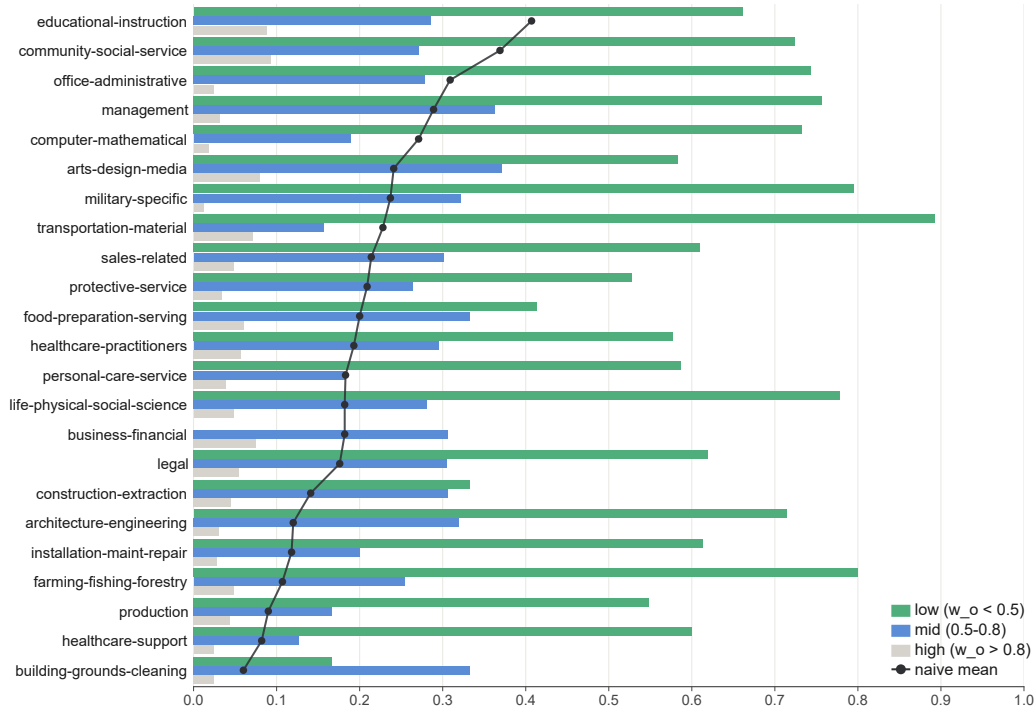
Configuration	n	wo avg	wo=1.0 %	PRG=0 (n)	PRG=0 breakdown (%)			Mean PRG
					Strict	Partial	LB no-gain	
Codex / GPT-5.4	653	0.753	37.7	296	83.1	15.5	1.4	0.200
Codex / GPT-5.1	676	0.608	30.0	300	67.7	18.7	13.7	0.288
Claude Code / Sonnet 4.6	584	0.776	49.8	317	91.8	7.3	0.9	0.185
Claude Code / Sonnet 4	653	0.636	23.6	245	62.9	29.0	8.2	0.272
OpenCode / GPT-5.4	631	0.695	30.7	281	69.0	19.6	11.4	0.213
OpenCode / Sonnet 4.6	644	0.771	45.8	343	86.0	8.2	5.8	0.171

395 **C.2 Stratified PRG by No-skill Baseline Strength**

396 We stratify Codex / GPT-5.4 scenarios by no-skill baseline pass rate into *low* ( $wo < 0.5$ ), *mid*  
 397 ( $0.5 \leq wo \leq 0.8$ ), and *high* ( $wo > 0.8$ ). Table A3 shows that the high stratum, which accounts for  
 398 58.8% of valid scenarios, dominates the unstratified mean; the unstratified value should therefore  
 399 not be read as the expected gain on challenging scenarios. Figure A2 reports the same diagnostic by  
 400 occupational category, where the relative category ranking is broadly preserved across strata even  
 401 though absolute magnitudes are compressed by ceiling effects.

**Table A3: Stratum-conditional PRG summary for Codex / GPT-5.4 ( $n = 653$  under Policy A).**

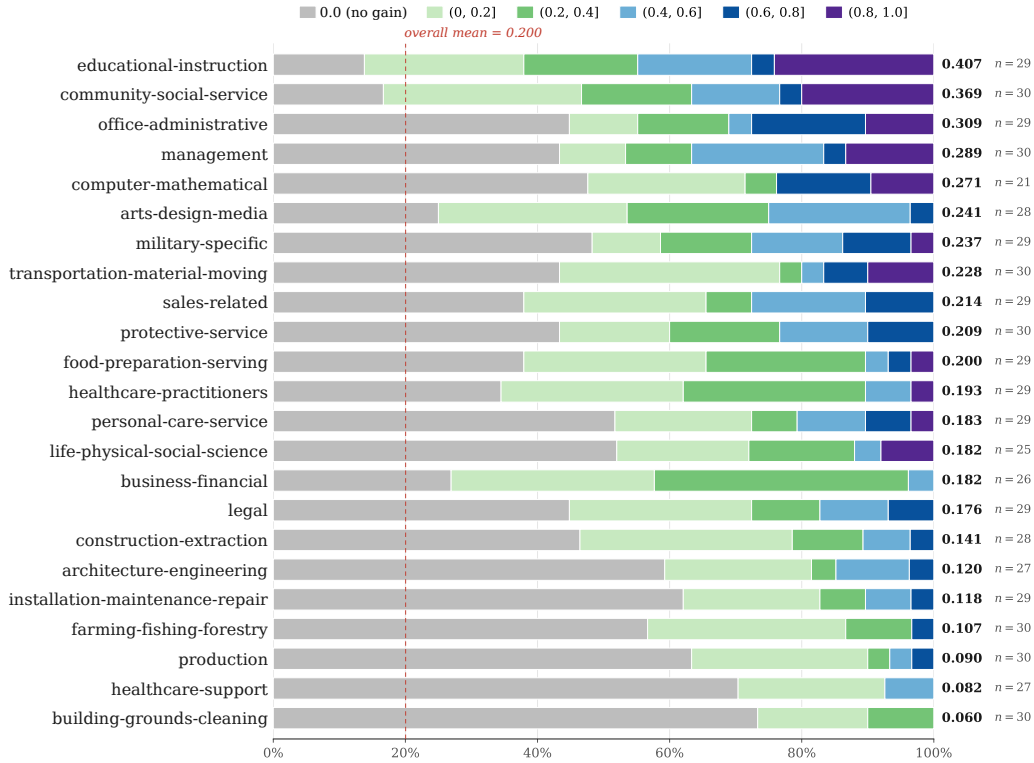
Stratum	n	mean PRG	median PRG	zero %
Low ( $wo < 0.5$ )	102	0.653	0.667	3.9
Mid ( $0.5 \leq wo \leq 0.8$ )	167	0.282	0.286	9.6
High ( $wo > 0.8$ )	384	0.045	0.000	71.9
All (unstratified)	653	0.200	0.143	45.3



**Figure A2: Stratified mean PRG by occupational category.** Rows report low-, mid-, and high-baseline strata for Codex / GPT-5.4; the black line marks the unstratified category mean.

402 **C.3 Per-category Distribution of PRG on a Single Configuration**

403 Figure A3 reports the full within-category distribution of scenario-level PRG for Codex / GPT-5.4,  
 404 complementing the mean heatmap in Figure 3 by exposing zero-gain mass and right-tail gains within  
 405 each category. The distribution separates three regimes that the category mean alone conceals:  
 406 *knowledge-structured* (high mean, small zero-gain share), *bimodal* (moderate mean coexisting with  
 407 a large zero-gain share and a strong right tail), and *operational or perception-driven* (low mean,  
 408 majority zero-gain). The zero-gain share is therefore a primary readout in its own right, separating  
 409 bimodal from knowledge-structured categories at comparable means.



**Figure A3: Per-scenario PRG distribution by occupational category.** Bars are stacked by gain bucket; right-side labels report category means. The dashed line marks the overall mean (0.200,  $n = 653$  under Policy A).

## 410 D Supplementary Safety Calibration

411 This appendix gives the calibration table used by the dynamic security judge when converting  
 412 behavioral evidence into exploitability scores.

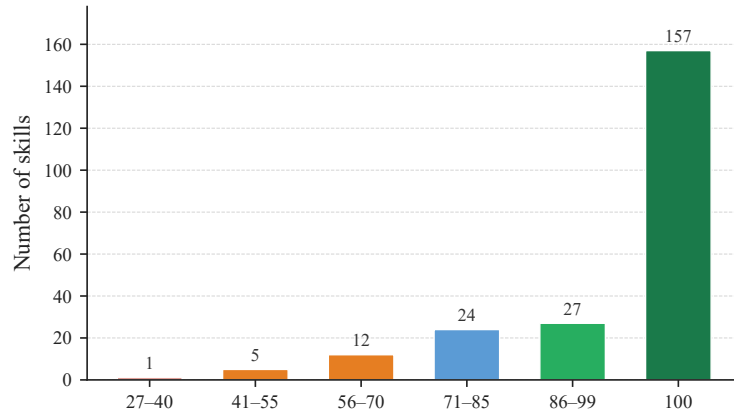
**Table A4:** Exploitability calibration for dynamic safety findings. The security judge maps each dynamic verdict to the exploitability range used in Eq. 5.

Verdict and score	Interpretation
confirmed (0.85–1.00)	End-to-end attack path is executed with observed evidence.
suspected (0.50–0.84)	Strong behavioral signal exists, but evidence is incomplete or ambiguous.
agent_refused (0.50)	The agent blocks the risky path, but the package-level risk remains.
path_exists_not_triggered (0.20–0.49)	The path appears reachable but is not exercised in this run.
likely_false_positive (0.10–0.19)	The trajectory shows little engagement with the predicted path.

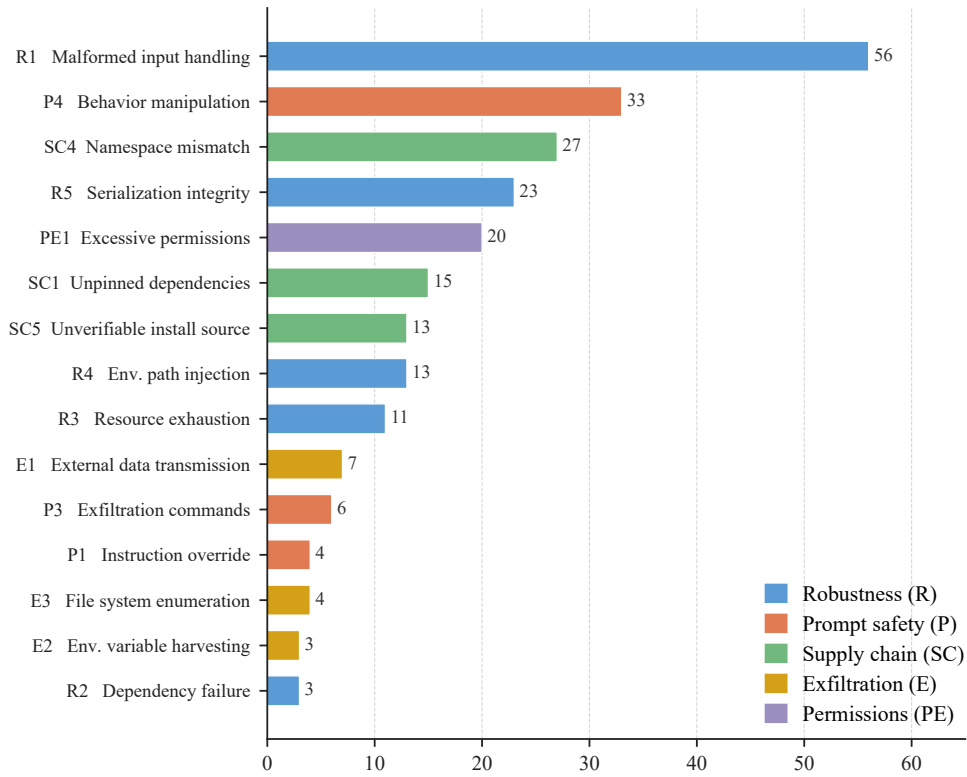
## 413 E Supplementary Safety Profile

414 This appendix complements the safety-profile analysis of Section 3.3. The main text reports agent-  
 415 conditioned dynamic outcomes in Figure 4; the diagnostics below extend the same analysis with  
 416 score-distribution, risk-pattern-, occupation-, and agent-resolved decompositions. Figure A4 reports  
 417 the skill-level safety-score distribution over the 226-skill set. Figure A5 reports the static-finding  
 418 frequency by risk pattern, and Figure A6 summarizes the aggregate occupational profile by separating  
 419 triggered risks, finding-only risks, and skills with no finding. Figure A7 further breaks down the  
 420 dynamic-probe triggered rate by occupational category and agent backbone, on the same fixed finding

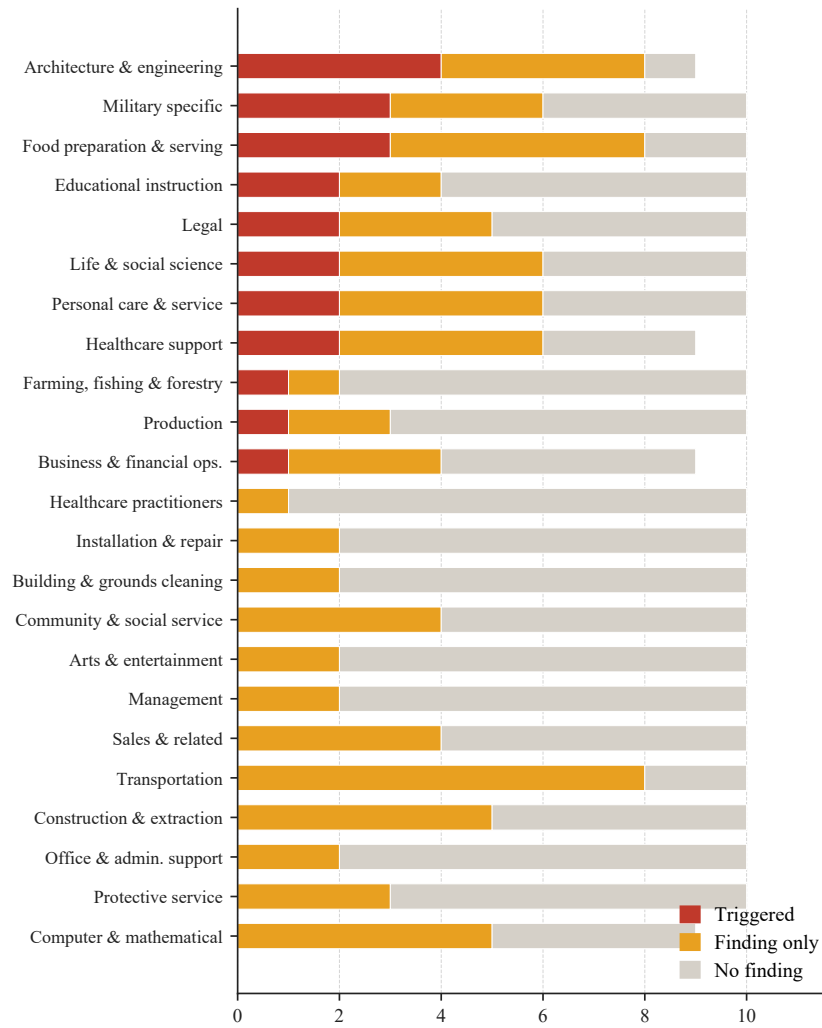
421 set used in Figure 4, so that cross-agent differences reflect agent behaviour rather than risk-set  
 422 differences.



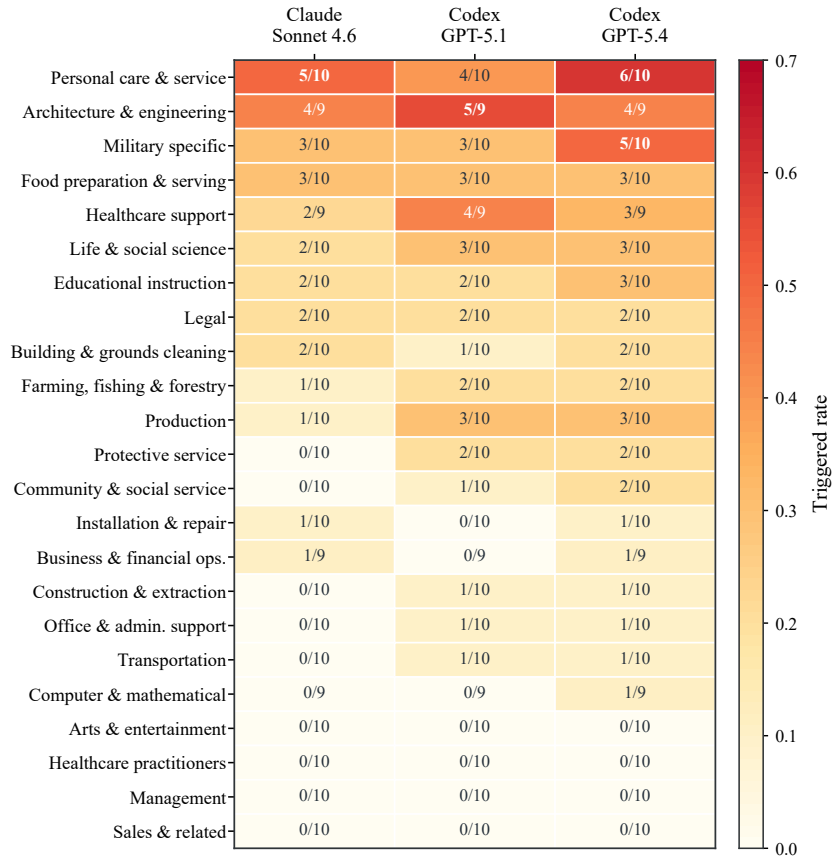
**Figure A4: Skill-level safety-score distribution.** Scores are computed over 226 evaluated skills. Most skills receive the maximum score, while the remaining skills form a risky tail.



**Figure A5: Static-scan finding frequency by risk pattern.** Bars report the number of static findings for each risk pattern over the 226-skill set, color-coded by risk family: Robustness (*R*), Prompt safety (*P*), Supply chain (*SC*), Exfiltration (*E*), and Permissions (*PE*).



**Figure A6: Risk profile by occupational category.** For each of the 23 occupational categories, skills are partitioned into *triggered* (a dynamic probe confirms exploitability), *finding only* (the scanner reports at least one static finding but no dynamic probe triggers), and *no finding*. Categories are sorted by the number of triggered skills.



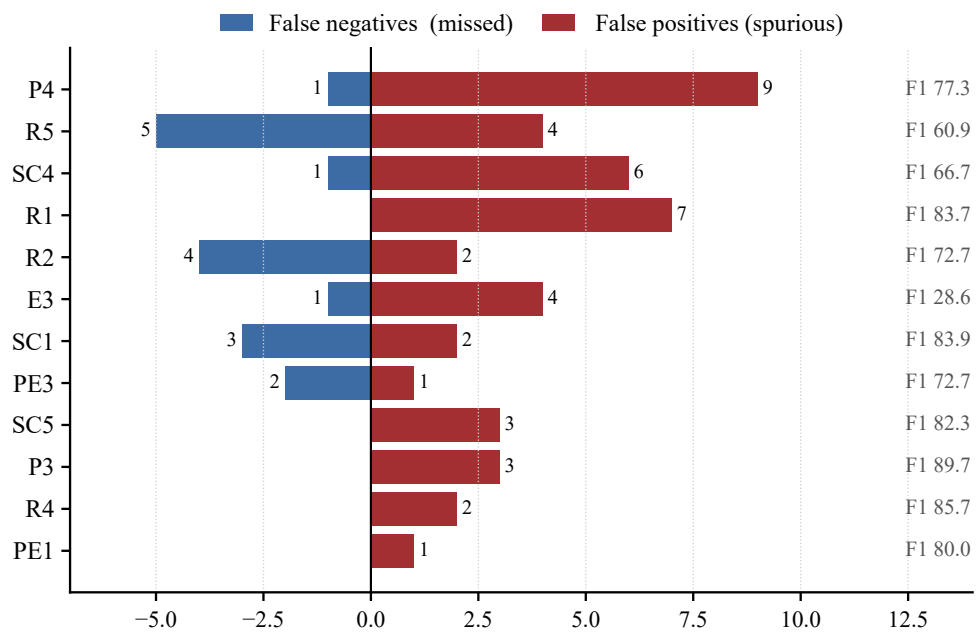
**Figure A7: Triggered rate by occupational category and agent.** Each cell reports triggered count over the number of tested skills in the category, on the same finding set as Figure 4(b). Rows are sorted by total triggered count across the three agents. Some categories remain consistently high-risk across agents, while others expose risk only under stronger execution backbones.

## 423 F Per-pattern Scanner Diagnostics

424 Table A5 reports the full per-pattern detection breakdown on the controlled subset of 29 skills. The  
 425 pattern-level error decomposition in Figure A8 highlights the asymmetry between over-reporting and  
 426 missed-detection error modes: P4, R1, and SC4 are dominated by false positives (over-reporting),  
 427 while R5, R2, and PE3 are dominated by false negatives (missed detections).

**Table A5:** Full per-pattern detection breakdown on the controlled subset of 29 skills.

ID	Pattern	Sev	#G	#P	TP	FP	FN	P (%)	R (%)	F1 (%)
<i>Prompt Injection</i>										
P1	Instruction Override	H	11	11	11	0	0	100.0	100.0	100.0
P2	Hidden Instructions	H	14	15	14	1	0	93.3	100.0	96.5
P3	Exfiltration Commands	H	13	16	13	3	0	81.2	100.0	89.7
P4	Behavior Manipulation	M	18	26	17	9	1	65.4	94.4	77.3
<i>Data Exfiltration</i>										
E1	External Data Transmission	M	16	17	16	1	0	94.1	100.0	97.0
E2	Env Var Harvesting	H	8	9	8	1	0	88.9	100.0	94.1
E3	FS Enumeration	M	2	5	1	4	1	20.0	50.0	28.6
E4	Context Leakage	H	3	4	3	1	0	75.0	100.0	85.7
<i>Privilege Escalation</i>										
PE1	Excessive Permissions	L	2	3	2	1	0	66.7	100.0	80.0
PE2	Sudo / Root Execution	M	1	1	1	0	0	100.0	100.0	100.0
PE3	Credential Access	H	6	5	4	1	2	80.0	66.7	72.7
<i>Supply Chain</i>										
SC1	Unpinned Dependencies	L	16	15	13	2	3	86.7	81.2	83.9
SC2	External Script Fetching	H	8	9	8	1	0	88.9	100.0	94.1
SC3	Obfuscated Code	H	3	3	3	0	0	100.0	100.0	100.0
SC4	Namespace Mismatch	L	8	13	7	6	1	53.8	87.5	66.7
SC5	Unverifiable Install Source	M	7	10	7	3	0	70.0	100.0	82.3
<i>Robustness</i>										
R1	Malformed Input Handling	L	18	25	18	7	0	72.0	100.0	83.7
R2	Network Timeout Missing	L	12	10	8	2	4	80.0	66.7	72.7
R3	Retry / Loop Bounds	L	2	3	2	1	0	66.7	100.0	80.0
R4	Error Handling	M	6	8	6	2	0	75.0	100.0	85.7
R5	Resource Cleanup	L	12	11	7	4	5	63.6	58.3	60.9



**Figure A8:** Per-pattern error decomposition for the 12 patterns with the largest total error (FP + FN). Negative bars indicate missed findings (false negatives); positive bars indicate spurious findings (false positives); pattern-level F1 is annotated on the right. The error profile is asymmetric: P4, R1, and SC4 are dominated by over-reporting, while R5, R2, and PE3 are dominated by missed detections.

## 428 **NeurIPS Paper Checklist**

429 The checklist is designed to encourage best practices for responsible machine learning research,  
430 addressing issues of reproducibility, transparency, research ethics, and societal impact. Do not remove  
431 the checklist: **The papers not including the checklist will be desk rejected.** The checklist should  
432 follow the references and follow the (optional) supplemental material. The checklist does NOT count  
433 towards the page limit.

434 Please read the checklist guidelines carefully for information on how to answer these questions. For  
435 each question in the checklist:

- 436 • You should answer [Yes], [No], or [N/A].
- 437 • [N/A] means either that the question is Not Applicable for that particular paper or the  
438 relevant information is Not Available.
- 439 • Please provide a short (1–2 sentence) justification right after your answer (even for [N/A]).

440 **The checklist answers are an integral part of your paper submission.** They are visible to the  
441 reviewers, area chairs, senior area chairs, and ethics reviewers. You will also be asked to include it  
442 (after eventual revisions) with the final version of your paper, and its final version will be published  
443 with the paper.

444 The reviewers of your paper will be asked to use the checklist as one of the factors in their evaluation.  
445 While [Yes] is generally preferable to [No], it is perfectly acceptable to answer [No] provided a  
446 proper justification is given (e.g., error bars are not reported because it would be too computationally  
447 expensive” or “we were unable to find the license for the dataset we used”). In general, answering  
448 [No] or [N/A] is not grounds for rejection. While the questions are phrased in a binary way, we  
449 acknowledge that the true answer is often more nuanced, so please just use your best judgment and  
450 write a justification to elaborate. All supporting evidence can appear either in the main paper or the  
451 supplemental material, provided in appendix. If you answer [Yes] to a question, in the justification  
452 please point to the section(s) where related material for the question can be found.

453 **IMPORTANT, please:**

- 454 • **Delete this instruction block, but keep the section heading “NeurIPS Paper Checklist”.**
- 455 • **Keep the checklist subsection headings, questions/answers and guidelines below.**
- 456 • **Do not modify the questions and only use the provided macros for your answers.**

### 457 **1. Claims**

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

460 Answer: **[TODO]**

461 Justification: **[TODO]**

462 Guidelines:

- 463 • The answer [N/A] means that the abstract and introduction do not include the claims  
464 made in the paper.
- 465 • The abstract and/or introduction should clearly state the claims made, including the  
466 contributions made in the paper and important assumptions and limitations. A [No] or  
467 [N/A] answer to this question will not be perceived well by the reviewers.
- 468 • The claims made should match theoretical and experimental results, and reflect how  
469 much the results can be expected to generalize to other settings.
- 470 • It is fine to include aspirational goals as motivation as long as it is clear that these goals  
471 are not attained by the paper.

### 472 **2. Limitations**

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

474 Answer: **[TODO]**

475 Justification: **[TODO]**

476  
477  
478  
479  
480  
481  
482  
483  
484  
485  
486  
487  
488  
489  
490  
491  
492  
493  
494  
495  
496  
497  
498  
499  
500  
501  
502  
503  
504  
505  
506  
507  
508  
509  
510  
511  
512  
513  
514  
515  
516  
517  
518  
519  
520  
521  
522  
523  
524  
525  
526

Guidelines:

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

3. **Theory assumptions and proofs**

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

Answer: **[TODO]**

Justification: **[TODO]**

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: **[TODO]**

Justification: **[TODO]**

Guidelines:

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

- 527
- 528
- 529
- 530
- 531
- 532
- 533
- 534
- 535
- 536
- 537
- 538
- 539
- 540
- 541
- 542
- 543
- 544
- 545
- 546
- 547
- 548
- 549
- 550
- 551
- 552
- 553
- 554
- 555
- 556
- 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.

## 557 5. Open access to data and code

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

561 Answer: [TODO]

562 Justification: [TODO]

563 Guidelines:

- 564
- 565
- 566
- 567
- 568
- 569
- 570
- 571
- 572
- 573
- 574
- 575
- 576
- 577
- 578
- 579
- 580
- 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).

- 581 • Providing as much information as possible in supplemental material (appended to the  
582 paper) is recommended, but including URLs to data and code is permitted.

## 583 6. Experimental setting/details

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

586 Answer: **[TODO]**

587 Justification: **[TODO]**

588 Guidelines:

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

## 594 7. Experiment statistical significance

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

597 Answer: **[TODO]**

598 Justification: **[TODO]**

599 Guidelines:

- 600 • The answer [N/A] means that the paper does not include experiments.
- 601 • The authors should answer **[Yes]** if the results are accompanied by error bars, confidence  
602 intervals, or statistical significance tests, at least for the experiments that support the  
603 main claims of the paper.
- 604 • The factors of variability that the error bars are capturing should be clearly stated (for  
605 example, train/test split, initialization, random drawing of some parameter, or overall  
606 run with given experimental conditions).
- 607 • The method for calculating the error bars should be explained (closed form formula,  
608 call to a library function, bootstrap, etc.)
- 609 • The assumptions made should be given (e.g., Normally distributed errors).
- 610 • It should be clear whether the error bar is the standard deviation or the standard error  
611 of the mean.
- 612 • It is OK to report 1-sigma error bars, but one should state it. The authors should  
613 preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis  
614 of Normality of errors is not verified.
- 615 • For asymmetric distributions, the authors should be careful not to show in tables or  
616 figures symmetric error bars that would yield results that are out of range (e.g., negative  
617 error rates).
- 618 • If error bars are reported in tables or plots, the authors should explain in the text how  
619 they were calculated and reference the corresponding figures or tables in the text.

## 620 8. Experiments compute resources

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

624 Answer: **[TODO]**

625 Justification: **[TODO]**

626 Guidelines:

- 627 • The answer [N/A] means that the paper does not include experiments.
- 628 • The paper should indicate the type of compute workers CPU or GPU, internal cluster,  
629 or cloud provider, including relevant memory and storage.
- 630 • The paper should provide the amount of compute required for each of the individual  
631 experimental runs as well as estimate the total compute.

- 632 • The paper should disclose whether the full research project required more compute  
633 than the experiments reported in the paper (e.g., preliminary or failed experiments that  
634 didn't make it into the paper).

635 **9. Code of ethics**

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

638 Answer: **[TODO]**

639 Justification: **[TODO]**

640 Guidelines:

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

647 **10. Broader impacts**

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

650 Answer: **[TODO]**

651 Justification: **[TODO]**

652 Guidelines:

- 653 • The answer [N/A] means that there is no societal impact of the work performed.
- 654 • If the authors answer [N/A] or [No], they should explain why their work has no societal  
655 impact or why the paper does not address societal impact.
- 656 • Examples of negative societal impacts include potential malicious or unintended uses  
657 (e.g., disinformation, generating fake profiles, surveillance), fairness considerations  
658 (e.g., deployment of technologies that could make decisions that unfairly impact specific  
659 groups), privacy considerations, and security considerations.
- 660 • The conference expects that many papers will be foundational research and not tied  
661 to particular applications, let alone deployments. However, if there is a direct path to  
662 any negative applications, the authors should point it out. For example, it is legitimate  
663 to point out that an improvement in the quality of generative models could be used to  
664 generate Deepfakes for disinformation. On the other hand, it is not needed to point out  
665 that a generic algorithm for optimizing neural networks could enable people to train  
666 models that generate Deepfakes faster.
- 667 • The authors should consider possible harms that could arise when the technology is  
668 being used as intended and functioning correctly, harms that could arise when the  
669 technology is being used as intended but gives incorrect results, and harms following  
670 from (intentional or unintentional) misuse of the technology.
- 671 • If there are negative societal impacts, the authors could also discuss possible mitigation  
672 strategies (e.g., gated release of models, providing defenses in addition to attacks,  
673 mechanisms for monitoring misuse, mechanisms to monitor how a system learns from  
674 feedback over time, improving the efficiency and accessibility of ML).

675 **11. Safeguards**

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

679 Answer: **[TODO]**

680 Justification: **[TODO]**

681 Guidelines:

- 682 • The answer [N/A] means that the paper poses no such risks.

- 683
- 684
- 685
- 686
- 687
- 688
- 689
- 690
- 691
- 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.

## 692 12. Licenses for existing assets

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

696 Answer: **[TODO]**

697 Justification: **[TODO]**

698 Guidelines:

- 699
- 700
- 701
- 702
- 703
- 704
- 705
- 706
- 707
- 708
- 709
- 710
- 711
- 712
- 713
- 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](https://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.

## 714 13. New assets

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

717 Answer: **[TODO]**

718 Justification: **[TODO]**

719 Guidelines:

- 720
- 721
- 722
- 723
- 724
- 725
- 726
- 727
- 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.

## 728 14. Crowdsourcing and research with human subjects

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

732 Answer: **[TODO]**

733 Justification: **[TODO]**

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

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?

Answer: **[TODO]**

Justification: **[TODO]**

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: **[TODO]**

Justification: **[TODO]**

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.