

---

# AgentPulse: A Continuous Multi-Signal Framework for Evaluating AI Agents in Deployment

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1       Static benchmarks measure what AI agents can do at a fixed point in time but  
2       not how they are adopted, maintained, or experienced in deployment. We intro-  
3       duce **AgentPulse**, a continuous evaluation framework scoring 50 agents across 10  
4       workload categories along four factors (Benchmark Performance, Adoption Sig-  
5       nals, Community Sentiment, and Ecosystem Health) aggregated from 18 real-time  
6       signals across GitHub, package registries, IDE marketplaces, social platforms, and  
7       benchmark leaderboards. Our headline finding, on  $n=34$  public-repo agents un-  
8       der a circularity-controlled test using only the Benchmark factor, is that bench-  
9       mark capability predicts whether an agent ships through public IDE marketplace  
10       extensions, replicated across three independent platforms (VS Code Marketplace,  
11       Open VSX, JetBrains Marketplace): Mann-Whitney  $U$  comparing  $B$ -scores of  
12       marketplace-present vs. marketplace-absent agents gives  $p=0.011, 0.0011, 0.004$   
13       with rank-biserial effect sizes  $+0.60, +0.86, +0.71$ . The split is between agents  
14       that publish public extensions and those that do not (a heterogeneous group span-  
15       ning closed-distribution agents, hosted services, and frameworks); the finding is  
16       therefore narrower than “commercialization” broadly construed. VS Code was  
17       pre-specified in the v6 draft; Open VSX and JetBrains are post-hoc external repli-  
18       cations. The same predictor is directionally negative on library-reuse (GitHub  
19       Dependents B-only  $\rho_s=-0.30$ , B+S  $\rho_s=-0.38$ ,  $p=0.04$ ), evidence of discrimi-  
20       native validity: the framework distinguishes “agent that completes tasks” from  
21       “library that gets reused.” AgentPulse is a methodology for surfacing deployment  
22       signal absent from benchmarks, not a ground-truth ranking. The framework, sig-  
23       nals, scoring outputs, and evaluation harness are released under CC BY 4.0.

## 24 1 Introduction

25       AI agents, systems that combine language model reasoning (17; 19; 20) with tool use (38) and multi-  
26       step planning (30), have moved from research prototypes to production tools used by millions of  
27       developers. Existing evaluation, however, remains anchored in static benchmarks: SWE-bench (6),  
28       GAIA (10), WebArena (16), AgentBench (8), and HumanEval+ (3) all measure capability at fixed  
29       points in time. These benchmarks are essential but incomplete: they cannot capture how agents  
30       evolve through frequent updates, how reliability varies under real-world load, or how developer  
31       experience changes as workflows adapt to these tools (91).

32       Consistent with the ED Track’s framing of evaluation as a scientific object of study, this paper  
33       investigates evaluation methodology for deployed AI agents. We ask: does a deployment-aware  
34       composite of public signals capture information benchmarks miss, and can such a composite be  
35       validated independently of the signals it aggregates? Our contribution is a measurement framework  
36       with three validation analyses plus a diagnostic, not a leaderboard.

37 We motivate this question through two concrete limitations of current agent evaluation:

38 Limitation 1: Existing benchmarks measure agents at fixed points along narrow capability slices,  
39 missing both the weekly release cadence of agentic systems and cross-tool comparability between,  
40 e.g., a coding copilot and a multi-agent framework whose workflows differ.

41 Limitation 2: Capability  $\neq$  adoption. An agent’s task-completion score does not predict which tools  
42 developers actually choose. Adoption depends on integration quality, pricing, reliability, documen-  
43 tation, and community, dimensions invisible to benchmarks.

44 **Empirical findings preview.** Three validation analyses plus one diagnostic ground the framework:  
45 (i) the four factors capture largely complementary information ( $n=50$ ;  $\rho_{\max}=0.75$  for Adoption-  
46 Ecosystem, all other pairwise  $|\rho|\leq 0.34$ ); (ii) within the  $n=34$  public-repo stratum, Benchmark ca-  
47 pability predicts the public-IDE-marketplace shipping decision across three independent platforms  
48 (Mann–Whitney  $U$  on  $B$ -scores for present vs. absent agents:  $p=0.011, 0.0011, 0.004$ ; effect sizes  
49  $r_{rb}=+0.60, +0.86, +0.71$ ); (iii) the same predictor is directionally negative on library-reuse met-  
50 rics, evidence of discriminative validity. The fourth analysis, a ranking-divergence diagnostic on the  
51  $n=11$  SWE-bench overlap (presented in §5.2 and dissected as a factor ablation in §6), is treated as  
52 a within-subset artifact rather than as a positive validity claim.

### 53 Contributions.

- 54 1. A four-factor evaluation framework, Benchmark Performance, Adoption Signals, Commu-  
55 nity Sentiment, and Ecosystem Health, with explicit design rationale (Section 3).
- 56 2. An 18-signal real-time data pipeline aggregating across GitHub, package registries, IDE  
57 marketplaces, social platforms (Bluesky, Reddit, Hacker News, Stack Overflow, Mastodon),  
58 and five benchmark leaderboards, evaluating 50 agents across 10 workload categories (Sec-  
59 tion 4).
- 60 3. Three validation analyses plus a diagnostic: factor independence ( $n=50$ ); circularity-  
61 controlled cross-factor predictive validity for the public-IDE-marketplace shipping deci-  
62 sion across three platforms ( $n=34$ ; Mann–Whitney  $U$   $p=0.011, 0.0011, 0.004$ ; rank-  
63 biserial  $r_{rb}=+0.60, +0.86, +0.71$ , all Holm-corrected); discriminative validity against  
64 library-reuse metrics ( $n=30$ ; B-only  $\rho_s=-0.30$ , B+S  $\rho_s=-0.38, p=0.04$ ); and a diagnos-  
65 tic ranking-divergence analysis on the  $n=11$  SWE-bench overlap (Section 5).
- 66 4. A factor ablation that transparently characterizes the composite’s behavior, including a  
67 structural tension on the  $n=11$  SWE-bench subset that we interpret rather than paper over  
68 (Section 6).
- 69 5. A public release of the framework, all collected signals, scored texts, and an evaluation  
70 harness designed for continuous, community-extensible agent evaluation (Appendix F).

## 71 2 Related Work

72 **Agent benchmarks.** SWE-bench (6) evaluates coding agents on real GitHub issue resolution;  
73 GAIA (10) tests general assistant capability across browsing, reasoning, and tool use; We-  
74 bArena (16) measures browser automation; AgentBench (8) provides a multi-environment evalua-  
75 tion suite; and HumanEval+ (3) extends code generation evaluation. These benchmarks measure  
76 what agents can do at a point in time, not how they are adopted and experienced over time. Agent-  
77 Pulse is complementary: it operates on top of these benchmarks as one of four factors and adds  
78 dimensions they do not capture.

79 **Holistic and continuous evaluation.** HELM (7) introduced multi-metric evaluation for language  
80 models. LMSYS Chatbot Arena (4; 15) pioneered continuous, preference-based LLM evaluation  
81 through pairwise human comparisons, providing a clear external ground truth that single-metric  
82 benchmarks lack. AgentPulse extends the holistic and continuous evaluation paradigm to agent  
83 ecosystems but takes a different validation approach: rather than soliciting human preferences, we  
84 ground the composite by testing that benchmark-and-sentiment factors alone predict observable  
85 adoption.

86 **Multi-signal aggregation.** Social-media sentiment (2), domain-specific NLP (1), and aspect-  
 87 based sentiment analysis (11) are established techniques in adjacent fields. We adapt these to agent  
 88 evaluation with agent-specific lexicons and combine them with adoption and ecosystem health met-  
 89 rics that have no analog in prior agent benchmarks. Our adoption signals draw conceptually on  
 90 classical models of technology diffusion (74) and on empirical software-engineering work mining  
 91 open-source ecosystems (79). Critiques of benchmark methodology (12) highlight precisely the  
 92 kinds of blind spots (adoption, deployment context, real-world signal) that AgentPulse is designed  
 93 to address.

### 94 3 The AgentPulse Framework

95 AgentPulse evaluates each agent through a composite of four factors:

$$AP(a) = w_B \cdot B(a) + w_A \cdot A(a) + w_S \cdot S(a) + w_E \cdot E(a) \quad (1)$$

96 where  $AP(a)$  denotes the composite (calligraphic) and  $S(a)$  the sentiment factor defined below, with  
 97  $w_B=0.35$ ,  $w_A=0.25$ ,  $w_S=0.20$ ,  $w_E=0.20$ . The allocation is a principled prior, not an empirical op-  
 98 timum: benchmark capability receives the largest weight because task completion is the most direct  
 99 measure of whether an agent does what it claims; the remaining 65% captures dimensions bench-  
 100 marks alone cannot measure. The framework supports custom weightings, and we test robustness  
 101 through sensitivity (Section 5.3) and ablation (Section 6).

102 **Factor 1: Benchmark Performance ( $B$ ).** The average of all available published benchmark  
 103 scores, normalized to  $[0, 1]$ :

$$B(a) = \frac{1}{|\mathcal{B}_a|} \sum_{b \in \mathcal{B}_a} \frac{s_b(a)}{100} \quad (2)$$

104 where  $\mathcal{B}_a$  is the set of benchmarks for which agent  $a$  has a published score among SWE-bench Veri-  
 105 fied, GAIA, WebArena, HumanEval+, and TAU-bench (14). Agents without published benchmarks  
 106 receive a neutral prior of 0.5 rather than the cross-sectional mean, ensuring the prior remains stable  
 107 as new agents enter the registry.

108 **Factor 2: Adoption Signals ( $A$ ).** Log-normalized metrics blending code hosting, package distri-  
 109 bution, and IDE penetration:

$$A(a) = 0.40 \cdot \hat{G}(a) + 0.35 \cdot \hat{D}(a) + 0.25 \cdot \hat{I}(a) \quad (3)$$

110 where  $\hat{G}(a) = \log_{10}(\text{stars}+1)/5.5$ ,  $\hat{D}(a) = \max(\log_{10}(D_{\text{pypi}}+1)/7, \log_{10}(D_{\text{npm}}+1)/7)$ , and  
 111  $\hat{I}(a) = \log_{10}(I_{\text{vsc}}+1)/8$ , with denominators corresponding to log-scale ceilings. Using  
 112  $\max(D_{\text{pypi}}, D_{\text{npm}})$  avoids penalizing agents distributed in only one ecosystem.

113 **Factor 3: Community Sentiment ( $S$ ).** Drawn from a multi-layer NLP pipeline (VADER (5),  
 114 TextBlob (9), FinBERT (1), and DistilBERT-SST2 (13)) applied to text from Bluesky, Reddit,  
 115 Hacker News, Stack Overflow, GitHub Discussions, Mastodon, Dev.to, V2EX, and Lemmy. Sentiment  
 116 is rescaled to  $[0, 1]$ :

$$S(a) = \text{clamp}(\bar{s}_{\text{composite}}(a) \cdot 2.5 + 0.5, 0, 1). \quad (4)$$

117 The multiplier 2.5 is calibrated to the empirical range of agent-level mean sentiment in our corpus:  
 118 although per-text sentiment (Appendix C) takes values on  $[-1, 1]$ , engagement-weighted means ag-  
 119 gregated across hundreds of mentions per agent fall within approximately  $[-0.2, 0.2]$ , so the affine  
 120 map sends typical values to the interior of  $[0, 1]$  without saturation. The pipeline applies sarcasm de-  
 121 tection, engagement weighting, and per-platform credibility weighting. Five domain-specific aspect  
 122 dimensions are detailed in Appendix C.

**Factor 4: Ecosystem Health ( $E$ ).**

$$E(a) = 0.3 \cdot C(a) + 0.2 \cdot r_{\text{close}} + 0.3 \cdot \max\left(1 - \frac{\Delta_{\text{days}}}{60}, 0\right) + 0.2 \cdot \frac{R_{\text{vsc}} - 2}{3} \quad (5)$$

123 where  $C(a) = \min(\log_{10}(\text{contributors}+1)/3, 1)$  is log-normalized contributor depth,  $r_{\text{close}}$  is the  
 124 GitHub issue close rate,  $\Delta_{\text{days}}$  is days since last update (60-day decay), and  $R_{\text{vsc}} \in [1, 5]$  is the  
 125 VS Code Marketplace rating.

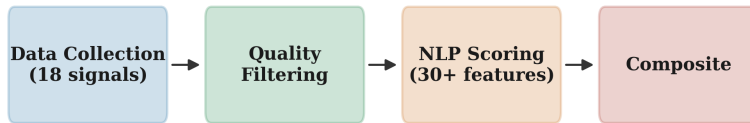


Figure 1: AgentPulse pipeline. Eighteen signals collected from public APIs and leaderboards are quality-filtered, scored through the NLP pipeline, and aggregated into the four-factor composite. The system runs autonomously and ingests new agents within hours.

Table 1: The 18 signals collected per agent, organized by factor. Sources are public APIs and leaderboard scrapes; collection cadence ranges from 5 minutes to 24 hours.

Factor	Signals	Source	Cadence
Benchmark (5)	SWE-bench, GAIA, WebArena, HumanEval+, TAU-bench	Leaderboards	24 hr
Adoption (6)	GitHub stars (+velocity), PyPI/npm downloads, VS Code installs/rating, Docker pulls	APIs	1 hr
Community (4)	Social sentiment, SO questions, issue close rate, contributors	APIs + NLP	1 hr
Ecosystem (3)	Days since release, doc depth proxy, enterprise-readiness composite	GitHub + MP	6 hr

126 **Design decisions.** No pricing factor: agents operate under heterogeneous pricing models, and  
 127 including pricing would structurally bias rankings toward open-source agents independent of ca-  
 128 pability. Closed-source measurement boundary: agents without public repositories or marketplace  
 129 presence receive zero on observable adoption sub-signals; we treat this as a measurement boundary  
 130 rather than a quality judgment (Section 7).

## 131 4 Data Pipeline and Agent Registry

132 We collect 18 signals per agent (Table 1); each collector runs independently with automatic retry.  
 133 Collected texts undergo MD5 and trigram-Jaccard deduplication, bot/spam filtering, and source-  
 134 credibility weighting before entering the NLP pipeline. The full data quality protocol is in Ap-  
 135 pendix A; the architecture is summarized in Figure 1. We track 50 agents across five functional  
 136 groups (development: 18; research & analysis: 6; browser: 7; multi-agent systems: 11; general: 8),  
 137 each mapped to one of 10 workload categories, enabling category-specific rankings (Appendix B).

138 **Registry composition and validation strata.** The 50-agent registry partitions naturally into three  
 139 subsets that the validation analyses use distinctly. The full  $n=50$  supports the factor-independence  
 140 analysis (Section 5.1), which is purely structural and uses no external target. The  $n=34$  public-repo  
 141 subset (agents with any observable GitHub repository in the v6 snapshot) is the appropriate scope  
 142 for the cross-factor predictive-validity test (Section 5.2), because all three IDE-marketplace proxies  
 143 are observable for this stratum (with zero-imputation when an agent has no marketplace presence).  
 144 The  $n=11$  subset of agents with published SWE-bench Verified scores is small and over-represents  
 145 closed-source/zero-adoption agents (Codex, Devin), so we treat it as a diagnostic only (Section 5.2,  
 146 Section 6); it is not the basis for any positive validity claim. We make this stratification explicit  
 147 because the  $n=11$  vs.  $n=34$  vs.  $n=50$  distinction does substantial conceptual work in what follows.

## 148 5 Validation

149 A multi-signal composite is justified only if (a) the factors capture genuinely complementary infor-  
 150 mation, and (b) the composite predicts something beyond the signals it aggregates. We test (a) at  
 151 full registry scale ( $n=50$ ), test (b) on the 34-agent public-repo subset using a circularity-controlled

Table 2: Inter-factor Spearman correlations across all 50 agents. Three factors are largely independent ( $|\rho| \leq 0.34$ ); only Adoption–Ecosystem co-vary substantially ( $\rho=0.75$ ). Adoption and Ecosystem both reflect open-source project health and are expected to correlate, but they are not redundant: Adoption measures demand (downloads, installs) while Ecosystem measures supply (contributors, maintenance).

	Benchmark ( $B$ )	Adoption ( $A$ )	Sentiment ( $S$ )	Ecosystem ( $E$ )
Benchmark ( $B$ )	1.00	0.19	0.04	0.34
Adoption ( $A$ )		1.00	0.10	<b>0.75</b>
Sentiment ( $S$ )			1.00	0.11
Ecosystem ( $E$ )				1.00

152 design, and finally examine ranking divergence on the  $n=11$  SWE-bench subset as a diagnostic of  
 153 the subset.

154 **What would convince us we’re wrong.** We commit in advance to three falsifiers: (i) if the IDE-  
 155 bucket shipping-decision effect failed to replicate at  $T+6$  on Open VSX or JetBrains, i.e., if the  
 156 Mann–Whitney  $p$ -value rose above 0.05 or the rank-biserial dropped below  $+0.2$ , the cross-platform  
 157 pre-specified claim would not survive; (ii) we predict that on a future snapshot with marketplace-  
 158 present subset size  $\geq 25$ , the within-presence rank correlation between  $B$  and install volume will  
 159 be positive and significant; if it is null or negative, the IDE result remains a discrete-presence effect  
 160 (a weaker but still meaningful claim) and our broader prediction about within-presence ranking is  
 161 falsified; (iii) if the library-reuse correlation flipped to positive, the discriminative-validity story  
 162 would be undermined. We commit to posting a held-out replication snapshot at month  $T+6$  post-  
 163 publication via the artifact repository, and the released code lets external readers run all three checks  
 164 directly on any snapshot.

## 165 5.1 Factor independence ( $n=50$ )

166 Pairwise Spearman correlations across all 50 agents (Table 2) show that the four factors carry largely  
 167 complementary signal. Benchmark–Sentiment ( $\rho=0.04$ ), Adoption–Sentiment ( $\rho=0.10$ ), and  
 168 Sentiment–Ecosystem ( $\rho=0.11$ ) are all near zero. Benchmark–Adoption ( $\rho=0.19$ ) and Benchmark–  
 169 Ecosystem ( $\rho=0.34$ ) are weak. The single notable correlation is Adoption–Ecosystem ( $\rho=0.75$ ),  
 170 which is expected on substantive grounds: both factors aggregate open-source project-health signals  
 171 from GitHub. They remain distinct sub-signals, Adoption captures user demand (stars, downloads,  
 172 installs) while Ecosystem captures maintainer supply (contributors, commit cadence, issue respon-  
 173 siveness), but they are not statistically independent and we treat them as a co-correlated pair rather  
 174 than as orthogonal factors.

## 175 5.2 Cross-factor predictive validity ( $n=34$ )

176 The most natural objection to a multi-signal composite is that it has no external ground truth. We ad-  
 177 dress this directly with a circularity-controlled design. Among the 50 agents, 34 have public GitHub  
 178 repositories and observable adoption signals in our snapshot. We use the *Benchmark factor alone* as  
 179 our primary circularity-controlled predictor (a maximally restrictive test:  $B$  contains no Adoption-,  
 180 Sentiment-, or Ecosystem-derived signals and has no causal path to platform adoption metrics) and  
 181 report the Benchmark+Sentiment (B+S) sub-composite as a robustness check on whether the broader  
 182 factor combination preserves the signal. We test whether each predicts external proxies grouped into  
 183 three theoretical buckets: **IDE-marketplace adoption** (developers actively installed the agent into  
 184 their IDE), **library-reuse adoption** (other repositories depend on the agent’s package), and **com-**  
 185 **munity engagement** (discussion volume on developer forums). All values are computed on the v6  
 186 release snapshot and are exactly reproducible via `python reproduce_table3.py`.

187 **IDE-marketplace shipping decision (the headline result).** Within the public-repo stratum,  
 188 benchmark capability predicts whether an agent ships an extension to public IDE marketplaces,  
 189 across three independent platforms. Mann–Whitney  $U$  on  $B$ -scores for marketplace-present vs.  
 190 absent agents (Table 3): VS Code  $U=42$ ,  $p=0.011$ ,  $r_{rb}=+0.60$ ; Open VSX  $U=12$ ,  $p=0.0011$ ,

Table 3: Cross-factor predictive validity ( $n=34$  public-repo agents). **IDE bucket primary test:** Mann–Whitney  $U$  of the Benchmark factor for marketplace-present vs. marketplace-absent agents (discrete public-marketplace shipping decision), reported as  $U / p / \text{rank-biserial } r_{rb}$ . The MW test does not apply outside the IDE bucket. **Secondary statistic:** Spearman  $\rho_s$  with B-only and B+S, 95% CIs from 1,000 bootstrap resamples (CI excluding zero implies  $p < 0.05$ ). Last column is non-zero marketplace-presence count out of  $n=34$  for IDE rows, and the actual evaluable subset size  $n$  for ‡ rows. VS Code was pre-specified in v6; Open VSX and JetBrains are post-hoc external replications. **Bold** marks the load-bearing argument.

Proxy	MW $U$ test ( $U / p / r_{rb}$ )	B-only $\rho_s$ [95% CI]	B+S $\rho_s$ [95% CI]	Pres./ $n$
<i>IDE-marketplace bucket</i>				
VS Code installs	<b>42 / 0.011 / +0.60</b>	+0.48 [+0.20,+0.70]	+0.42 [+0.14,+0.65]	8
Open VSX installs	<b>12 / 0.0011 / +0.86</b>	+0.60 [+0.35,+0.79]	+0.55 [+0.29,+0.72]	6
JetBrains downloads	<b>27 / 0.004 / +0.71</b>	+0.52 [+0.28,+0.72]	+0.50 [+0.26,+0.69]	7
<i>Library-reuse bucket (negative = discriminative validity)</i>				
GitHub Dependents ‡	—	−0.30 [−0.61,+0.04]	<b>−0.38 [−0.68,−0.01]</b>	30
PyPI downloads/month ‡	—	+0.09 [−0.41,+0.61]	−0.10 [−0.60,+0.47]	14
<i>Community-engagement bucket (null)</i>				
GitHub stars ‡	—	0.00 [−0.42,+0.41]	−0.03 [−0.44,+0.37]	31
Stack Overflow questions (raw)	—	+0.29 [−0.08,+0.56]	+0.18 [−0.19,+0.50]	34
Discord member count	—	+0.11 [−0.20,+0.42]	−0.02 [−0.35,+0.29]	34

191  $r_{rb}=+0.86$ ; JetBrains  $U=27$ ,  $p=0.004$ ,  $r_{rb}=+0.71$ . VS Code was pre-specified in the v6 draft;  
192 Open VSX (used by VSCodium and JetBrains Fleet) and JetBrains are post-hoc external replica-  
193 tions, all three significant under Holm correction. Spearman  $\rho_s=0.48$ – $0.60$  (B-only) and  $0.42$ – $0.55$   
194 (B+S) are reported as confirmatory secondary statistics. *Scope of the claim.* The marketplace-absent  
195 group is heterogeneous; it includes agents distributed via closed APIs, hosted-only services, and pure  
196 frameworks. The finding is therefore best read as “among public-repo agents, those with stronger  
197 benchmarks are more likely to ship public IDE extensions,” not as a general commercialization-  
198 decision claim. *Mechanism.* The Benchmark factor uses only published leaderboard scores and  
199 has no mechanical link to marketplace presence; cross-platform replication on three independent  
200 ecosystems is the strongest evidence against alternative explanations.

201 **Discriminative validity against library-reuse metrics.** GitHub Dependents (count of public  
202 repositories importing the agent’s package) is negatively correlated with the predictors:  $\rho_s=-0.30$   
203 for B alone ( $p=0.11$ ,  $n=30$ , directionally consistent but underpowered) and  $\rho_s=-0.38$  for B+S  
204 ( $p=0.04$ ). Library-reuse metrics are dominated by frameworks (LangGraph 38,792 dependents,  
205 Claude MCP 50,709, CrewAI 18,319), which lack published benchmarks and receive the 0.5 prior,  
206 while benchmark-strong standalone agents (Claude Code, Tabnine) are not consumed as libraries.  
207 The negative direction is the expected discriminative-validity outcome: AgentPulse correctly distin-  
208 guishes “agent that completes tasks” from “library that gets reused.”

209 **Community-engagement nulls.** GitHub stars, Stack Overflow question volume, and Discord  
210 member counts all fail to replicate (Table 3). We attribute these to the same registry-composition  
211 effect that drives the library-reuse negatives plus, for SO specifically, tag-alias noise (the cursor  
212 SO tag also matches CSS-cursor questions). Appendix E reports per-proxy Pearson, leave-one-out,  
213 and alternative-sub-composite robustness checks.

214 **What this analysis does not show.** (a) The test is restricted to the  $n=34$  public-repo stratum;  
215 closed-source agents are unobservable in our measurement scope and the result does not extend  
216 to them. (b) Within the marketplace-present subset ( $n=6$ – $8$ ), the rank correlation between B and  
217 install volume is non-significant; we therefore claim only the discrete shipping decision, not within-  
218 presence install-volume ranking. (c) Temporal stability is untested in this snapshot; we commit to a  
219 held-out replication at  $T+6$  months. Further caveats appear in Appendix E.

220 **Robustness to weight perturbation.** We resample the weight vector 1,000 times from  
221  $\text{Dir}(3.5, 2.5, 2.0, 2.0)$ , centered on our default allocation, concentration  $\kappa=10$ , and recompute the  
222 B+S Spearman against each IDE proxy. All 1,000 draws are positive for all three proxies. Leave-  
223 one-out across 34 single-agent drops yields B-only  $\rho_s$  ranges  $[0.42, 0.55]$  (VS Code),  $[0.55, 0.65]$

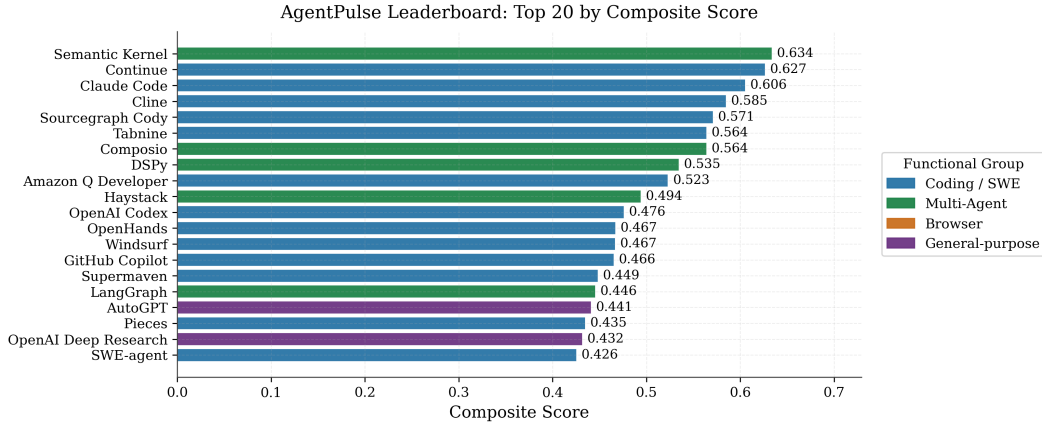


Figure 2: AgentPulse leaderboard: top 20 agents by composite score across the 50-agent registry, colored by functional group. SWE agents, open-source coding agents, browser agents, multi-agent frameworks, copilots, and general-purpose agents all appear in the top 20, with the dominant differentiating factor varying by type. Closed-API agents appear lower because their adoption signals are unobservable under our measurement scope, not because their capability is lower (Section 7).

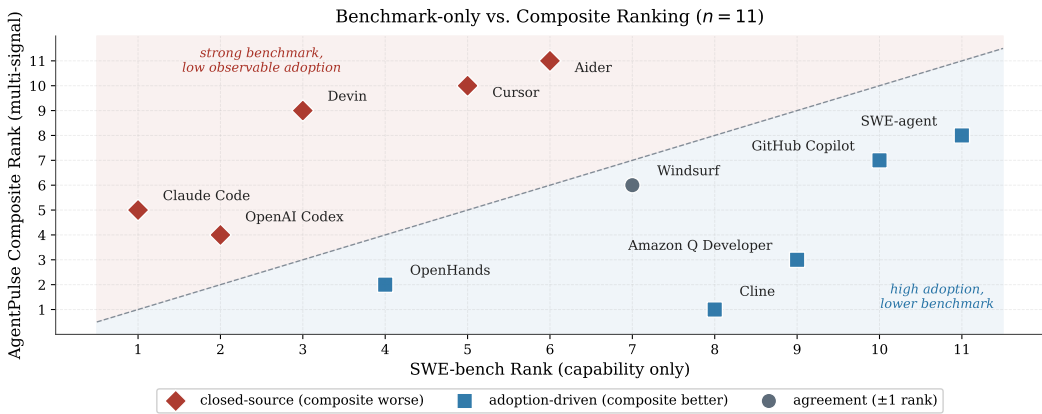


Figure 3: Benchmark-only vs. composite ranking for the 11 agents with published SWE-bench scores. Blue points (Cline, OpenHands, Amazon Q Developer) rank higher on the composite than on benchmarks alone, lifted by adoption and ecosystem signal. Red points (Aider, Cursor, Devin) rank lower on the composite, primarily because they expose limited observable adoption signal under our measurement scope (closed APIs and proprietary tools). Grey points indicate agreement within  $\pm 1$  rank. The Spearman correlation between the two orderings on this subset is  $\rho_s=0.09$ .

224 (Open VSX), [0.46, 0.58] (JetBrains), with  $|\Delta| \leq 0.07$  per proxy. Neither the weight allocation nor  
 225 any single agent drives the IDE-bucket signal.

226 **Ranking divergence on the SWE-bench overlap (diagnostic).** We finally examine ranking diver-  
 227 gence on the  $n=11$  subset of agents with published SWE-bench Verified scores (Figure 3). Among  
 228 these 11 agents, benchmark-only and composite rankings disagree on 25 of 55 pairwise compar-  
 229 isons, with 7 of 11 agents shifting by  $\geq 2$  rank positions. The low Spearman correlation between  
 230 composite and benchmark-only orderings ( $\rho_s=0.09$ ) reflects the composite’s deliberate incorpora-  
 231 tion of non-benchmark dimensions. We frame this analysis as diagnostic given the sample size; the  
 232 evidence for the framework’s value is Section 5.2. The same closed-source/open-source asymmetry  
 233 that drives the  $n=11$  ablation pattern in Section 6 also drives the divergence statistics here, which is  
 234 why we treat Section 5.2 as the framework’s primary validity claim.

235 The divergences cluster in two regimes:

Table 4: Top agents by composite score in three SE-relevant categories. Different categories surface different leaders, reflecting that the framework distinguishes capability profiles across distinct workflows. **Bold** marks the per-category leader.

Category	Agent	Composite	$B$	$A$	$S$
Coding	<b>Claude Code</b>	<b>0.602</b>	0.824	0.368	0.580
	Cline	0.585	0.565	0.553	0.545
	OpenHands	0.530	0.615	0.273	0.883
SWE	<b>Claude Code</b>	<b>0.602</b>	0.824	0.368	0.580
	OpenAI Codex	0.464	0.796	0.000	0.578
Multi-Agent	<b>OpenAI Agents SDK</b>	<b>0.577</b>	0.500	0.321	0.560
	LangGraph	0.573	0.500	0.444	0.500

- 236 • Cline ranks 7th on SWE-bench (38.0%) but 2nd on composite within this subset, driven  
237 by 3.7M VS Code Marketplace installs ( $A=0.553$ , the highest among coding agents) and  
238 200+ contributors. This places weight on tool-preference factors beyond autonomous task  
239 completion.
- 240 • OpenAI Codex ranks 2nd on SWE-bench (69.1%) but 4th on composite. As a closed API  
241 without a public repository, package, or extension, it has no observable adoption signal, a  
242 measurement boundary rather than a quality verdict.
- 243 • Devin ranks 3rd on SWE-bench (55.0%) but 9th on composite, for the same closed-  
244 distribution reason.

245 **Category-specific rankings.** The framework yields different leaders in different workload categories  
246 (Table 4), consistent with the cross-tool comparability motivation. Within-category rankings  
247 sometimes invert overall composite ranking, evidence the framework captures category-relevant signal  
248 rather than a single global ordering.

### 249 5.3 Sensitivity to weight perturbations

Table 5: Single-factor sensitivity: rank change ( $\Delta$ ) for five representative agents when each factor weight is increased by +10pp and the other three are reduced proportionally. Claude Code’s leading position in the SWE category is invariant; no agent shifts by more than one rank.

Agent	$B \uparrow$	$A \uparrow$	$S \uparrow$	$E \uparrow$
<b>Claude Code</b>	<b>0</b>	<b>0</b>	<b>0</b>	<b>0</b>
Cline	-1	+1	0	+1
OpenHands	+1	-1	+1	-1
GitHub Copilot	0	0	0	+1
OpenAI Codex	0	-1	0	-1

250 We perturb each factor weight by  $\pm 10$  percentage points, redistributing proportionally (Table 5).  
251 The SWE-category leader is invariant across all perturbations, and no agent shifts by more than  
252 one rank position under any single-factor perturbation. Bootstrap confidence intervals on per-agent  
253 composite scores (1,000 resamples of the underlying signal data) show no agent in the top 20 shifts  
254 in median composite by more than  $\pm 0.018$ .

## 255 6 Diagnostic: Why the SWE-bench Subset Misleads

256 The pathology of the  $n=11$  SWE-bench-scored subset is itself a motivation for a multi-signal frame-  
257 work, not a failure of one. We make this argument concretely. The full composite is nearly uncorre-  
258 lated with SWE-bench within this subset ( $\rho=0.03$ ), while benchmark-only ranking achieves  $\rho=0.95$   
259 (Table 6). This is striking but not paradoxical: the  $n=11$  subset over-represents closed-source/high-  
260 capability/zero-adoption agents (Codex, Devin) opposite high-adoption/moderate-benchmark ones  
261 (Cline, Copilot), producing a within-subset negative Adoption–Capability correlation that does not

### Agent Factor Decomposition — Top 12 Agents

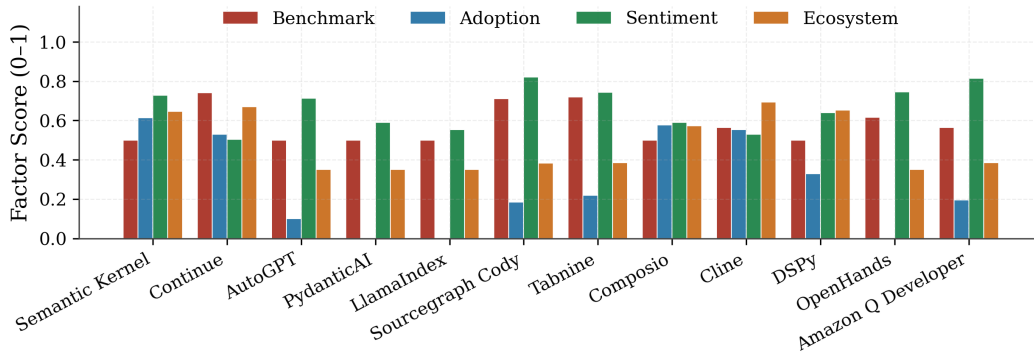


Figure 4: Factor decomposition for the top 12 agents. Different agents are differentiated by different factors: SWE-agent and OpenHands are lifted by sentiment (green); coding-focused agents like Claude Code by benchmark performance (red); copilots and multi-agent frameworks by adoption (blue). This decomposition motivates the multi-factor composite over a single score.

Table 6: Factor ablation ( $n=11$  SWE-bench-scored agents; over-represents closed-source/zero-adoption agents and is diagnostic only, the headline validation is Table 3,  $n=34$ ). Each row zeroes one factor, redistributes its weight proportionally, and recomputes  $\rho_s$  vs. SWE-bench Verified.

Scheme	$w_B$	$w_A$	$w_S$	$w_E$	$\rho_s$ vs. SWE-bench
Full composite	0.35	0.25	0.20	0.20	<b>0.03</b>
w/o Benchmark	0	0.38	0.31	0.31	-0.33
w/o Adoption	0.47	0	0.27	0.27	0.57
w/o Sentiment	0.44	0.31	0	0.25	0.42
w/o Ecosystem	0.44	0.31	0.25	0	0.10
Benchmark-only	1.00	0	0	0	<b>0.95</b>

262 generalize; the cross-factor analysis on  $n=34$  (Section 5.2) shows Benchmark capability *positively*  
 263 predicts IDE marketplace presence across the broader registry. Substantively, removing Benchmark  
 264 ( $\rho=-0.33$ ) produces rankings anti-correlated with capability (Benchmark is non-substitutable),  
 265 while Adoption captures genuinely orthogonal signal; Cline’s 3.7M VS Code installs despite moder-  
 266 ate SWE-bench is paradigmatic. The ablation table is therefore a diagnostic for why benchmark-only  
 267 evaluation is insufficient on heterogeneous registries, not a verdict on the full composite.

## 268 7 Discussion and Limitations

269 **What AgentPulse is and is not.** AgentPulse is a measurement methodology with circularity-  
 270 controlled validation that it captures externally observable adoption information beyond benchmarks  
 271 alone. It is not a developer-preference oracle, a substitute for capability benchmarks, or an evaluator  
 272 of closed-source agents whose adoption is unobservable; the IDE-bucket test uses the  $n=34$  public-  
 273 repo subset to avoid this confound, and rankings of closed-source agents should be read alongside  
 274 benchmark-only rankings. Downstream users may reweight factors via the released harness.

275 **Sentiment retention, limitations, and what’s next.** Sentiment’s aggregate Spearman with IDE in-  
 276 stalls is near zero ( $|\rho_s|<0.13$ ); its retained value is at the aspect level (Appendix C, Table 9: Devin’s  
 277 reliability-aspect sentiment is negative despite a moderate aggregate).  $w_S=0.20$  may be too high for  
 278 capability-focused use cases. Other limitations: posters are not representative, calibration ( $\kappa=0.81$ ,  
 279  $n=200$ ) is a sanity check, the pipeline is English-only, and only 11 agents have SWE-bench scores.  
 280 Beyond the pre-registered  $T+6$  replication, planned next steps include a within-presence rank test  
 281 once the marketplace-present subset reaches  $n\geq 25$  and a developer-preference survey for validation  
 282 independent of all 18 signals. AgentPulse is released under CC BY 4.0 (App. F).

283 **References**

- 284 [1] D. Araci. FinBERT: Financial Sentiment Analysis with Pre-Trained Language Models. *arXiv*  
285 *preprint arXiv:1908.10063*, 2019.
- 286 [2] J. Bollen, H. Mao, and X. Zeng. Twitter mood predicts the stock market. *Journal of Computa-*  
287 *tional Science*, 2(1):1–8, 2011.
- 288 [3] M. Chen et al. Evaluating Large Language Models Trained on Code. *arXiv preprint*  
289 *arXiv:2107.03374*, 2021.
- 290 [4] W.-L. Chiang et al. Chatbot Arena: An Open Platform for Evaluating LLMs by Human Prefer-  
291 ence. *ICML*, 2024.
- 292 [5] C. J. Hutto and E. Gilbert. VADER: A Parsimonious Rule-based Model for Sentiment Analysis  
293 of Social Media Text. *ICWSM*, 2014.
- 294 [6] C. E. Jimenez et al. SWE-bench: Can Language Models Resolve Real-World GitHub Issues?  
295 *ICLR*, 2024.
- 296 [7] P. Liang et al. Holistic Evaluation of Language Models. *Annals of the New York Academy of*  
297 *Sciences*, 2023.
- 298 [8] X. Liu et al. AgentBench: Evaluating LLMs as Agents. *ICLR*, 2024.
- 299 [9] S. Loria. TextBlob: Simplified Text Processing. <https://textblob.readthedocs.io>,  
300 2018.
- 301 [10] G. Mialon et al. GAIA: A Benchmark for General AI Assistants. *arXiv preprint*  
302 *arXiv:2311.12983*, 2023.
- 303 [11] M. Pontiki et al. SemEval-2016 Task 5: Aspect Based Sentiment Analysis. *SemEval*, 2016.
- 304 [12] I. D. Raji, E. M. Bender, et al. AI and the Everything in the Whole Wide World Benchmark.  
305 *NeurIPS*, 2021.
- 306 [13] V. Sanh, L. Debut, J. Chaumond, and T. Wolf. DistilBERT, a Distilled Version of BERT:  
307 Smaller, Faster, Cheaper and Lighter. *arXiv preprint arXiv:1910.01108*, 2019.
- 308 [14] S. Yao et al.  $\tau$ -bench: A Benchmark for Tool-Agent-User Interaction in Real-World Domains.  
309 *arXiv preprint arXiv:2406.12045*, 2024.
- 310 [15] L. Zheng et al. Judging LLM-as-a-Judge with MT-Bench and Chatbot Arena. *NeurIPS*, 2023.
- 311 [16] S. Zhou et al. WebArena: A Realistic Web Environment for Building Autonomous Agents.  
312 *ICLR*, 2024.
- 313 [17] A. Vaswani et al. Attention Is All You Need. *NeurIPS*, 2017.
- 314 [18] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova. BERT: Pre-training of Deep Bidirectional  
315 Transformers for Language Understanding. *NAACL-HLT*, 2019.
- 316 [19] T. Brown et al. Language Models are Few-Shot Learners. *NeurIPS*, 2020.
- 317 [20] J. Achiam et al. GPT-4 Technical Report. *arXiv preprint arXiv:2303.08774*, 2023.
- 318 [21] H. Touvron et al. LLaMA: Open and Efficient Foundation Language Models. *arXiv preprint*  
319 *arXiv:2302.13971*, 2023.
- 320 [22] H. Touvron et al. Llama 2: Open Foundation and Fine-Tuned Chat Models. *arXiv preprint*  
321 *arXiv:2307.09288*, 2023.
- 322 [23] Y. Liu et al. RoBERTa: A Robustly Optimized BERT Pretraining Approach. *arXiv preprint*  
323 *arXiv:1907.11692*, 2019.
- 324 [24] C. Raffel et al. Exploring the Limits of Transfer Learning with a Unified Text-to-Text Trans-  
325 former. *Journal of Machine Learning Research*, 21(140):1–67, 2020.

- 326 [25] M. Lewis et al. BART: Denoising Sequence-to-Sequence Pre-training for Natural Language  
327 Generation, Translation, and Comprehension. *ACL*, 2020.
- 328 [26] P. He, X. Liu, J. Gao, and W. Chen. DeBERTa: Decoding-Enhanced BERT with Disentangled  
329 Attention. *ICLR*, 2021.
- 330 [27] K. Clark, M.-T. Luong, Q. V. Le, and C. D. Manning. ELECTRA: Pre-training Text Encoders  
331 as Discriminators Rather than Generators. *ICLR*, 2020.
- 332 [28] Z. Lan et al. ALBERT: A Lite BERT for Self-supervised Learning of Language Representa-  
333 tions. *ICLR*, 2020.
- 334 [29] N. Reimers and I. Gurevych. Sentence-BERT: Sentence Embeddings using Siamese BERT-  
335 Networks. *EMNLP-IJCNLP*, 2019.
- 336 [30] J. Wei et al. Chain-of-Thought Prompting Elicits Reasoning in Large Language Models.  
337 *NeurIPS*, 2022.
- 338 [31] J. Wei et al. Emergent Abilities of Large Language Models. *Transactions on Machine Learning  
339 Research*, 2022.
- 340 [32] J. Kaplan et al. Scaling Laws for Neural Language Models. *arXiv preprint arXiv:2001.08361*,  
341 2020.
- 342 [33] J. Hoffmann et al. Training Compute-Optimal Large Language Models. *NeurIPS*, 2022.
- 343 [34] L. Ouyang et al. Training Language Models to Follow Instructions with Human Feedback.  
344 *NeurIPS*, 2022.
- 345 [35] P. Christiano, J. Leike, T. Brown, M. Martic, S. Legg, and D. Amodei. Deep Reinforcement  
346 Learning from Human Preferences. *NeurIPS*, 2017.
- 347 [36] Y. Bai et al. Constitutional AI: Harmlessness from AI Feedback. *arXiv preprint  
348 arXiv:2212.08073*, 2022.
- 349 [37] S. Yao et al. ReAct: Synergizing Reasoning and Acting in Language Models. *ICLR*, 2023.
- 350 [38] T. Schick et al. Toolformer: Language Models Can Teach Themselves to Use Tools. *NeurIPS*,  
351 2023.
- 352 [39] G. Wang, Y. Xie, Y. Jiang, A. Mandlekar, C. Xiao, Y. Zhu, L. Fan, and A. Anandkumar. Voy-  
353 ager: An Open-Ended Embodied Agent with Large Language Models. *Transactions on Ma-  
354 chine Learning Research*, 2024.
- 355 [40] J. S. Park, J. C. O’Brien, C. J. Cai, M. R. Morris, P. Liang, and M. S. Bernstein. Generative  
356 Agents: Interactive Simulacra of Human Behavior. *UIST*, 2023.
- 357 [41] N. Shinn, F. Cassano, A. Gopinath, K. Narasimhan, and S. Yao. Reflexion: Language Agents  
358 with Verbal Reinforcement Learning. *NeurIPS*, 2023.
- 359 [42] O. Khattab et al. DSPy: Compiling Declarative Language Model Calls into Self-Improving  
360 Pipelines. *ICLR*, 2024.
- 361 [43] P. Lewis et al. Retrieval-Augmented Generation for Knowledge-Intensive NLP Tasks. *NeurIPS*,  
362 2020.
- 363 [44] V. Karpukhin et al. Dense Passage Retrieval for Open-Domain Question Answering. *EMNLP*,  
364 2020.
- 365 [45] D. Hendrycks et al. Measuring Massive Multitask Language Understanding. *ICLR*, 2021.
- 366 [46] A. Srivastava et al. Beyond the Imitation Game: Quantifying and Extrapolating the Capabilities  
367 of Language Models. *Transactions on Machine Learning Research*, 2023.
- 368 [47] A. Wang et al. SuperGLUE: A Stickier Benchmark for General-Purpose Language Understand-  
369 ing Systems. *NeurIPS*, 2019.

- 370 [48] A. Wang, A. Singh, J. Michael, F. Hill, O. Levy, and S. R. Bowman. GLUE: A Multi-Task  
371 Benchmark and Analysis Platform for Natural Language Understanding. *ICLR*, 2019.
- 372 [49] P. Rajpurkar, J. Zhang, K. Lopyrev, and P. Liang. SQuAD: 100,000+ Questions for Machine  
373 Comprehension of Text. *EMNLP*, 2016.
- 374 [50] K. Cobbe et al. Training Verifiers to Solve Math Word Problems. *arXiv preprint*  
375 *arXiv:2110.14168*, 2021.
- 376 [51] D. Hendrycks et al. Measuring Mathematical Problem Solving with the MATH Dataset.  
377 *NeurIPS Datasets and Benchmarks*, 2021.
- 378 [52] J. Austin et al. Program Synthesis with Large Language Models. *arXiv preprint*  
379 *arXiv:2108.07732*, 2021.
- 380 [53] R. Li et al. StarCoder: May the Source Be with You! *Transactions on Machine Learning*  
381 *Research*, 2023.
- 382 [54] B. Rozière et al. Code Llama: Open Foundation Models for Code. *arXiv preprint*  
383 *arXiv:2308.12950*, 2023.
- 384 [55] Z. Chen et al. T-Eval: Evaluating the Tool Utilization Capability of Large Language Models  
385 Step by Step. *ACL*, 2024.
- 386 [56] Y. Qin et al. ToolLLM: Facilitating Large Language Models to Master 16000+ Real-World  
387 APIs. *ICLR*, 2024.
- 388 [57] S. G. Patil, T. Zhang, X. Wang, and J. E. Gonzalez. Gorilla: Large Language Model Connected  
389 with Massive APIs. *NeurIPS*, 2024.
- 390 [58] R. Bommasani et al. On the Opportunities and Risks of Foundation Models. *arXiv preprint*  
391 *arXiv:2108.07258*, 2021.
- 392 [59] E. M. Bender, T. Gebru, A. McMillan-Major, and S. Shmitchell. On the Dangers of Stochastic  
393 Parrots: Can Language Models Be Too Big? *FAccT*, 2021.
- 394 [60] M. Mitchell et al. Model Cards for Model Reporting. *FAccT*, 2019.
- 395 [61] T. Gebru et al. Datasheets for Datasets. *Communications of the ACM*, 64(12):86–92, 2021.
- 396 [62] L. Weidinger et al. Ethical and Social Risks of Harm from Language Models. *arXiv preprint*  
397 *arXiv:2112.04359*, 2021.
- 398 [63] B. Pang and L. Lee. Opinion Mining and Sentiment Analysis. *Foundations and Trends in In-*  
399 *formation Retrieval*, 2(1–2):1–135, 2008.
- 400 [64] B. Liu. Sentiment Analysis and Opinion Mining. *Synthesis Lectures on Human Language Tech-*  
401 *nologies*, 5(1):1–167, 2012.
- 402 [65] S. M. Mohammad and P. D. Turney. Crowdsourcing a Word–Emotion Association Lexicon.  
403 *Computational Intelligence*, 29(3):436–465, 2013.
- 404 [66] R. Socher et al. Recursive Deep Models for Semantic Compositionality Over a Sentiment Tree-  
405 bank. *EMNLP*, 2013.
- 406 [67] A. L. Maas, R. E. Daly, P. T. Pham, D. Huang, A. Y. Ng, and C. Potts. Learning Word Vectors  
407 for Sentiment Analysis. *ACL*, 2011.
- 408 [68] T. Mikolov, K. Chen, G. Corrado, and J. Dean. Efficient Estimation of Word Representations  
409 in Vector Space. *ICLR Workshop*, 2013.
- 410 [69] J. Pennington, R. Socher, and C. D. Manning. GloVe: Global Vectors for Word Representation.  
411 *EMNLP*, 2014.
- 412 [70] M. Honnibal and I. Montani. spaCy 2: Natural Language Understanding with Bloom Embed-  
413 dings, Convolutional Neural Networks and Incremental Parsing. Software, 2017.

- 414 [71] S. Bird, E. Klein, and E. Loper. *Natural Language Processing with Python*. O’Reilly Media,  
415 2009.
- 416 [72] A. Joulin, E. Grave, P. Bojanowski, and T. Mikolov. Bag of Tricks for Efficient Text Classifica-  
417 tion. *EACL*, 2017.
- 418 [73] M. Grootendorst. BERTopic: Neural Topic Modeling with a Class-Based TF-IDF Procedure.  
419 *arXiv preprint arXiv:2203.05794*, 2022.
- 420 [74] E. M. Rogers. *Diffusion of Innovations*. Free Press, 5th edition, 2003.
- 421 [75] F. D. Davis. Perceived Usefulness, Perceived Ease of Use, and User Acceptance of Information  
422 Technology. *MIS Quarterly*, 13(3):319–340, 1989.
- 423 [76] F. M. Bass. A New Product Growth for Model Consumer Durables. *Management Science*,  
424 15(5):215–227, 1969.
- 425 [77] V. Venkatesh, M. G. Morris, G. B. Davis, and F. D. Davis. User Acceptance of Information  
426 Technology: Toward a Unified View. *MIS Quarterly*, 27(3):425–478, 2003.
- 427 [78] A. Mockus, R. T. Fielding, and J. D. Herbsleb. Two Case Studies of Open Source Software  
428 Development: Apache and Mozilla. *ACM Transactions on Software Engineering and Methodol-  
429 ogy*, 11(3):309–346, 2002.
- 430 [79] E. Kalliamvakou, G. Gousios, K. Blincoe, L. Singer, D. M. German, and D. Damian. The  
431 Promises and Perils of Mining GitHub. *Mining Software Repositories (MSR)*, 2014.
- 432 [80] K. Crowston, K. Wei, J. Howison, and A. Wiggins. Free/Libre Open-Source Software Devel-  
433 opment: What We Know and What We Do Not Know. *ACM Computing Surveys*, 44(2):1–35,  
434 2012.
- 435 [81] B. Vasilescu, Y. Yu, H. Wang, P. Devanbu, and V. Filkov. Quality and Productivity Outcomes  
436 Relating to Continuous Integration in GitHub. *ESEC/FSE*, 2015.
- 437 [82] H. Borges, A. Hora, and M. T. Valente. Understanding the Factors that Impact the Popularity  
438 of GitHub Repositories. *ICSME*, 2016.
- 439 [83] J. Howard and S. Ruder. Universal Language Model Fine-tuning for Text Classification. *ACL*,  
440 2018.
- 441 [84] M. E. Peters et al. Deep Contextualized Word Representations. *NAACL-HLT*, 2018.
- 442 [85] A. Radford, J. Wu, R. Child, D. Luan, D. Amodei, and I. Sutskever. Language Models are  
443 Unsupervised Multitask Learners. *OpenAI Technical Report*, 2019.
- 444 [86] T. Wolf et al. Transformers: State-of-the-Art Natural Language Processing. *EMNLP System  
445 Demonstrations*, 2020.
- 446 [87] A. Paszke et al. PyTorch: An Imperative Style, High-Performance Deep Learning Library.  
447 *NeurIPS*, 2019.
- 448 [88] Y. Liu et al. A Survey on Evaluation of Large Language Models. *ACM Transactions on Intelli-  
449 gent Systems and Technology*, 15(3):1–45, 2024.
- 450 [89] Y. Chang et al. A Survey on Evaluation of Large Language Models. *ACM Transactions on  
451 Intelligent Systems and Technology*, 15(3):1–45, 2024.
- 452 [90] T. Guo et al. Large Language Model Based Multi-Agents: A Survey of Progress and Chal-  
453 lenges. *IJCAI*, 2024.
- 454 [91] Z. Xi et al. The Rise and Potential of Large Language Model Based Agents: A Survey. *Science  
455 China Information Sciences*, 2025.

## 456 Organization of Appendix

457 The appendix is organized as follows: in App. A we document the data-quality protocol applied to every  
458 collected text before NLP scoring (composite quality score, per-platform statistics, edge cases);  
459 in App. B we list the full 50-agent registry with workload-category mappings, signal-source assignments,  
460 and registry-construction criteria; in App. C we detail the NLP pipeline (sentiment composite,  
461 aspect dimensions, sarcasm detection, engagement weighting); in App. D we report the full sensitivity  
462 analyses (single- and multi-factor weight perturbations, bootstrap confidence intervals, robustness  
463 to subset selection); in App. E we expand the cross-factor predictive validity methodology, why  
464 the analysis matters, how the Benchmark+Sentiment sub-composite is constructed, how the external  
465 proxies were chosen, full results with robustness checks, and what the analysis does not show; in  
466 App. F we describe the released artifact (release-artifact structure, reproduction commands, compute  
467 requirements, license, anonymization for review); the NeurIPS Paper Checklist is appended last.

## 468 A Data Quality Protocol

469 This appendix documents the data-quality layer applied to every collected text before it enters the  
470 NLP scoring pipeline (Section 3). The goal is to ensure that downstream sentiment and aspect  
471 scores reflect substantive developer discussion rather than spam, duplicate boilerplate, bot-generated  
472 content, or topically irrelevant text.

### 473 A.1 Composite quality score

474 Four sub-scores are computed per text and combined into a single quality score:

$$q(t) = \theta_{\text{uniq}} q_{\text{uniq}}(t) + \theta_{\text{bot}} q_{\text{bot}}(t) + \theta_{\text{cred}} q_{\text{cred}}(t) + \theta_{\text{spec}} q_{\text{spec}}(t) \quad (6)$$

475 with  $\theta_{\text{uniq}} = \theta_{\text{spec}} = 0.30$ ,  $\theta_{\text{bot}} = \theta_{\text{cred}} = 0.20$ , and exclusion threshold  $q^* = 0.30$ . Texts with  $q(t) < q^*$  are  
476 excluded from sentiment scoring but retained in the released artifact (with a quality flag) for reviewer  
477 inspection.

478 **Uniqueness ( $q_{\text{uniq}}$ ).** A two-stage check. Stage 1 (exact): MD5 hash on lowercased, URL-stripped,  
479 whitespace-normalized text. Exact duplicates score 0. Stage 2 (near-duplicate): Trigram-Jaccard  
480 similarity against all previously seen texts within a 7-day rolling window, with threshold  $\tau = 0.85$ .  
481 The score is  $1 - J(t, t')$  where  $J(t, t')$  is the maximum Jaccard similarity to any prior text  $t'$ ; if  
482 no near-duplicate exists, the score is 1.0. The 7-day rolling window balances duplicate detection  
483 against memory cost; longer windows produce diminishing returns on flagging rate.

484 **Bot detection ( $q_{\text{bot}}$ ).** Heuristic score combining four signals: (i) content length (texts under  
485 20 characters or over 5,000 characters score lower); (ii) match against a curated list of  $\sim 40$   
486 spam/promotional regex patterns (e.g., affiliate-link patterns, repeated-emoji patterns, “DM me to  
487 learn more” patterns); (iii) author posting frequency relative to platform median (authors posting  
488 more than  $10 \times$  median rate flagged); (iv) engagement anomalies (texts with implausibly low or high  
489 engagement relative to author history). The score is the arithmetic mean of the four sub-signals,  
490 each in  $[0, 1]$ .

491 **Source credibility ( $q_{\text{cred}}$ ).** Per-platform base credibility weight reflecting moderation rigor and  
492 signal-to-noise, multiplied by a post-level engagement booster. Base weights are: Stack Over-  
493 flow 0.90, GitHub Discussions 0.85, Hacker News 0.85, Reddit 0.70, Mastodon 0.65, Lemmy 0.60,  
494 Bluesky 0.60, Dev.to 0.55, V2EX 0.60. The booster is  $\min(1.0, 0.5 + 0.1 \log_{10}(\text{engagement} + 1))$ ,  
495 so a post with 100 upvotes/likes receives a 0.7 booster, a post with 1,000 receives 0.8, etc. Final  
496 credibility is base  $\times$  booster, clamped to  $[0, 1]$ .

497 **Specificity ( $q_{\text{spec}}$ ).** Match rate against a curated list of  $\sim 180$  technical terms covering five cat-  
498 egories: (i) model names (Claude, GPT-4, Llama, etc.); (ii) framework names (LangGraph, Lla-  
499 maIndex, AutoGen, etc.); (iii) benchmark names (SWE-bench, GAIA, WebArena, etc.); (iv) integer  
500 pricing patterns (e.g.,  $\backslash\$\d{+}/\text{month}$ ); (v) version-number patterns (e.g.,  $\backslash\d{+}\.\d{+}$ ). The score is  
501  $\min(\text{matches}/3, 1.0)$ , so a text with three or more technical-term matches receives full credit.

502 **A.2 Per-platform statistics**

503 Across 15,000 collected texts, 39 (0.26%) were flagged for exclusion. Per-platform statis-  
504 tics appear in Table 7. The flagged texts all matched the multi-criterion pattern [duplicate,  
505 bot\_suspected, too\_generic], indicating that exclusions are concentrated in obviously low-  
506 signal content rather than borderline cases.

Table 7: Data-quality statistics by source platform. Higher quality and specificity indicate more reliable signal. Lower-moderation platforms (Dev.to, V2EX) flag more frequently.

Platform	$n$	Quality $\bar{q}$	Uniqueness	Bot score	Specificity	Flagged %
Bluesky	6,525	0.622	0.804	0.588	0.480	0.12
Hacker News	4,390	0.569	0.449	0.584	0.376	0.00
Reddit	1,217	0.669	0.982	0.579	0.396	0.08
arXiv	1,019	0.475	0.395	0.599	0.410	0.10
Dev.to	589	0.396	0.037	0.571	0.336	4.58
Stack Overflow	567	0.664	0.716	0.596	0.377	0.00
Mastodon	357	0.601	0.597	0.600	0.545	0.00
GitHub Disc.	258	0.626	0.667	0.566	0.358	0.00
V2EX	75	0.539	0.547	0.530	0.337	2.67
Lemmy	3	0.706	1.000	0.600	0.467	0.00
All	15,000	0.605	0.659	0.586	0.452	0.26

507 **A.3 Edge cases and failure modes**

508 We document three known failure modes of the quality pipeline:

509 (i) Cross-lingual content. Non-English text matching English technical terms by chance receives  
510 elevated specificity scores. Currently mitigated by language detection on text  $>50$  characters; texts  
511 identified as non-English are excluded from sentiment scoring (their counts contribute to engage-  
512 ment statistics only).

513 (ii) Long-form blog posts. Posts with thousands of words and only a few mentions of the agent of  
514 interest may receive high specificity scores due to dense technical vocabulary in unrelated sections.  
515 Mitigated by computing specificity over a  $\pm 200$ -word window centered on the agent mention rather  
516 than the entire text.

517 (iii) Adversarial templated praise. Coordinated promotional posts that vary surface text but share  
518 substantive content evade exact-duplicate detection. The trigram-Jaccard threshold catches most  
519 such cases ( $\tau=0.85$  flags posts sharing  $\geq 85\%$  of trigrams); we acknowledge this is an arms race  
520 and note it as a limitation.

521 **B Agent Registry and Workload Categories**

522 This appendix documents the agent registry, the inclusion criteria, and the 10 workload categories  
523 used for category-specific rankings.

524 **B.1 Registry construction**

525 The registry was constructed via combined automated discovery and manual curation. Automated  
526 discovery scans the OpenRouter catalog and GitHub-trending repositories every 6 hours; new can-  
527 didates are filtered against the inclusion criteria below before being added. Manual curation verifies  
528 provider attribution, primary workload category, and signal-source mappings.

529 **Inclusion criteria.** Agents were included if they:

- 530 1. were publicly available (i.e., usable by an external developer, whether free or paid);
- 531 2. had at least one observable signal among the 18 (a published benchmark, public repository,  
532 package distribution, marketplace listing, or social-platform mention);

533 3. primarily targeted agentic workflows, defined as multi-step task completion involving tool  
534 use, code execution, or autonomous decision-making, rather than chat-only interaction.

535 **Excluded categories.** Three categories were explicitly excluded:

- 536 • Superseded model versions (e.g., GPT-3.5 once GPT-4 was released; we track only the  
537 current production version per provider).
- 538 • Free community variants of paid agents (to avoid double-counting, e.g., we track Cursor  
539 but not derived community forks).
- 540 • Specialized non-text models (image generators, audio agents, video agents) outside the  
541 scope of agentic SE/general workflows.

## 542 **B.2 Workload categories**

543 Each agent is mapped to one of 10 workload categories. The categories are:

- 544 • coding, IDE-integrated coding assistants (autocomplete, refactor, inline edit) for individual  
545 files.
- 546 • copilot, copilot-style assistants tightly integrated with a development workflow, typically  
547 with chat + edit modes.
- 548 • swe, autonomous software engineering agents that resolve issues end-to-end (read repo,  
549 plan changes, edit, test).
- 550 • multi, multi-agent orchestration frameworks (build agents that coordinate sub-agents).
- 551 • browser, agents that operate a browser to complete web tasks (form filling, scraping, navi-  
552 gation).
- 553 • research, research-and-summarization agents (long-form synthesis from multiple web  
554 sources).
- 555 • enterprise, agents targeting enterprise integrations (knowledge base search, internal tool  
556 orchestration).
- 557 • general, general-purpose autonomous task agents without a specific workflow specializa-  
558 tion.
- 559 • consumer, conversational consumer assistants (broad audience, not agentic specialization).
- 560 • data, data-analysis-focused agents (notebooks, structured data Q&A, charting).

## 561 **B.3 Full registry**

## 562 **B.4 Signal availability**

563 Of the 50 agents:

- 564 • 11 have published SWE-bench Verified scores (used for the  $n=11$  ranking-divergence anal-  
565 ysis in Section 5.2): Claude Code, OpenAI Codex, Devin, OpenHands, Cursor, Windsurf,  
566 Cline, GitHub Copilot, SWE-agent, Aider, Amazon Q Developer.
- 567 • 35 have any observable GitHub repository (used for cross-factor predictive validity in Sec-  
568 tion 5.2).
- 569 • 16 have repositories with  $\geq 1,000$  stars (the typical threshold above which open-source  
570 signal becomes noise-resistant).
- 571 • 11 are distributed via the VS Code Marketplace, including Cline, Cursor, GitHub Copilot,  
572 Continue, Tabnine, Sourcegraph Cody, Supermaven, and others (full list in the released  
573 registry).

574 Closed-source agents (Devin, OpenAI Codex, Cursor, OpenAI Deep Research, Operator, Manus)  
575 lack observable signals in the Adoption factor specifically (no public repository, no package distri-  
576 bution, no marketplace listing); they may still register in Sentiment and partial Ecosystem signals.  
577 ChatGPT, for example, has no Adoption-factor footprint but very high Sentiment-factor mention  
578 volume, which is reflected in its overall composite ranking. This is the measurement boundary  
579 discussed in Section 7.

Table 8: Full 50-agent registry with provider attribution and primary workload category.

Group	Agent	Provider	Primary category
Development (18)	Claude Code	Anthropic	swe
	Cursor	Anysphere	coding
	OpenAI Codex	OpenAI	coding
	GitHub Copilot	GitHub	copilot
	Windsurf	Codeium	copilot
	Gemini CLI	Google	copilot
	Cline	Cline	coding
	Devin	Cognition	swe
	Replit Agent	Replit	coding
	OpenHands	OpenHands	coding
	SWE-agent	Princeton	swe
	Aider	Aider	coding
	Bolt	StackBlitz	coding
	Continue	Continue	coding
	Amazon Q Developer	Amazon	coding
	Tabnine	Tabnine	coding
	Sourcegraph Cody	Sourcegraph	coding
	Supermaven	Supermaven	coding
Research & Analysis (6)	OpenAI Deep Research	OpenAI	enterprise
	Perplexity Research	Perplexity	research
	Gemini Deep Research	Google	research
	Manus	Manus AI	general
	NotebookLM	Google	research
	Genspark	Genspark	research
Browser (7)	OpenClaw	Anthropic	browser
	Operator	OpenAI	browser
	Wingman	Wingman	general
	Browser Use	Browser Use	browser
	Adept ACT-2	Adept	browser
	NanoBot	NanoBot	browser
	Multion	Multion	browser
Multi-Agent Systems (11)	LangGraph	LangChain	multi
	CrewAI	CrewAI	enterprise
	Microsoft AutoGen	Microsoft	enterprise
	OpenAI Agents SDK	OpenAI	multi
	Claude MCP	Anthropic	multi
	Semantic Kernel	Microsoft	enterprise
	LlamaIndex	LlamaIndex	multi
	PydanticAI	Pydantic	multi
	DSPy	Stanford	multi
	Haystack	deepset	multi
	Composio	Composio	multi
General (8)	ChatGPT	OpenAI	consumer
	Claude	Anthropic	data
	AutoGPT	AutoGPT	general
	MetaGPT	MetaGPT	general
	Lovable	Lovable	coding
	v0	Vercel	coding
	Pieces	Pieces	coding
	Kimi Researcher	Moonshot	research

580 **C NLP Pipeline Detail**

581 This appendix documents the full NLP pipeline, including the sentiment composite, calibration  
 582 methodology, aspect dimensions, and post-processing.

Agent Scores by Category — Top Agents per Workload

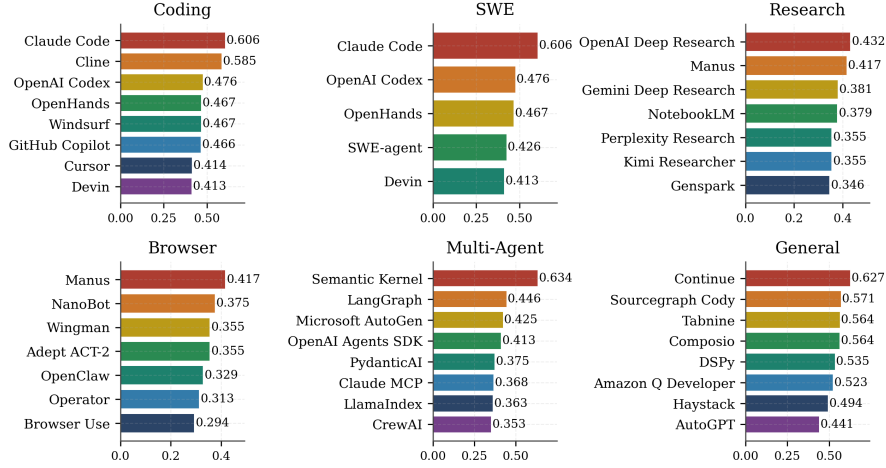


Figure 5: Per-category top agents. The framework yields different leaders in different workload categories: Claude Code dominates SWE; multi-agent frameworks (LangGraph, PydanticAI) lead the multi-agent category; OpenClaw and NanoBot lead the browser category. Within-category rankings sometimes invert overall composite ranking, evidence the framework captures category-relevant signal rather than a single global ordering.

### 583 C.1 Sentiment composite

584 The per-text composite sentiment is

$$S_i = 0.40 s_{\text{VADER}} + 0.20 s_{\text{TB}} + 0.25 s_{\text{FinBERT}} + 0.15 s_{\text{DistilBERT}}. \quad (7)$$

585 The four components capture complementary aspects of polarity:

- 586 • VADER (5) (weight 0.40): rule-based, fast, calibrated for social-media text, handles emoji  
587 and intensifiers well.
- 588 • TextBlob (9) (weight 0.20): pattern-based, handles formal text robustly, complements  
589 VADER on long-form posts.
- 590 • FinBERT (1) (weight 0.25): fine-tuned for evaluative-financial text; particularly useful for  
591 pricing discussions and adoption-cost framing.
- 592 • DistilBERT-SST2 (13) (weight 0.15): general-purpose neural sentiment, captures patterns  
593 the lexicon-based methods miss.

594 **Calibration.** The blending weights were calibrated on a held-out validation set of 200 manu-  
595 ally labeled LLM/agent discussion posts drawn proportionally from each platform. The composite  
596 achieves 89% agreement with human polarity labels (vs. 82% for VADER alone, 85% for FinBERT  
597 alone). Inter-annotator agreement is Cohen’s  $\kappa=0.81$  (two annotators; disagreements resolved by  
598 discussion).

599 **Limitation.** The 200-text calibration set is small and was annotated by the authors. We treat  
600 this as an internal sanity check on pipeline behavior, not as a benchmark of the sentiment factor’s  
601 external accuracy. Larger external annotation efforts would strengthen the validation; we welcome  
602 replication.

### 603 C.2 Aspect dimensions

604 For each of five LLM/agent-specific aspects  $a \in \{\text{performance, reliability, cost, innovation, adoption}\}$ ,  
605 we maintain positive lexicon  $L_a^+$  ( $\sim 30\text{--}60$  terms each) and negative lexicon  $L_a^-$  of similar size.

606 Per-text aspect score is

$$s_a(t) = \frac{n_a^+(t) - n_a^-(t)}{n_a^+(t) + n_a^-(t)}, \quad I_a(t) = \min\left(\frac{n_a^+(t) + n_a^-(t)}{\delta_a}, 1\right) \quad (8)$$

607 where  $n_a^\pm(t)$  counts term occurrences in text  $t$ , and  $\delta_a$  is an intensity-saturation constant tuned  
 608 per-aspect ( $\delta_{\text{performance}}=5$ , others  $\delta=3$ ). Both raw scores and intensities are released for all 15,000  
 609 NLP-scored texts.

610 **Lexicon construction.** Lexicons were seeded from a hand-curated set of evaluative terms drawn  
 611 from publicly available developer discussion (blog posts, README files, benchmark-paper related-  
 612 work sections), expanded via word2vec nearest-neighbor expansion (cosine similarity threshold  
 613 0.65) within the pre-collected developer-text corpus, then manually filtered by both annotators. Each  
 614 lexicon was reviewed for face validity; ambiguous terms (e.g., “fast” for performance vs. “fast-  
 615 shipping” for adoption) were assigned to the most-frequent contextual sense. The final lexicons and  
 616 their construction notebook are released as part of the artifact for inspection and modification.

617 **Aspect-level findings.** Aggregate sentiment masks dimension-specific strengths and weaknesses  
 618 (Table 9). Claude Code leads on code quality and debugging; Cline dominates IDE integration (con-  
 619 sistent with its VS Code adoption); Devin shows negative reliability sentiment, suggesting agentic-  
 620 loop instability invisible in its moderate aggregate score.

Table 9: Aspect-level sentiment for three coding agents. Devin’s negative reliability score is masked by its moderate aggregate sentiment, illustrating why aspect decomposition matters. Devin’s IDE-integration cell is dashed because the aspect lexicon yielded fewer than the  $n=30$ -mention threshold required to report a stable per-aspect mean (Devin operates as a hosted agent rather than an IDE extension).

Aspect	Claude Code	Cline	Devin
Code quality	+0.18	+0.14	+0.11
Debugging	+0.12	+0.08	-0.03
Multi-file editing	+0.09	+0.06	+0.02
Agentic reliability	+0.07	+0.10	-0.12
IDE integration	+0.05	+0.22	n/a (insufficient mentions)

### 621 C.3 Sarcasm detection

622 A list of  $K=24$  regex patterns capturing ironic praise, explicit markers (r"/s\b"), and contradictory  
 623 polarity-content combinations. Patterns include:

- 624 • Ironic praise patterns (e.g., r"(amazing|brilliant|genius)\b.\*\b(not|never|nope)");
- 625 • Explicit sarcasm markers (r"/s\b", "/sarcasm", "/j/k");
- 626 • Negated superlatives (e.g., “definitely not the best...”);
- 627 • Excessive intensifiers (e.g., “soooo great”), with manual review of false-positive rate.

628 Texts with  $P_{\text{sarcasm}}(t) > 0.30$  have their sentiment sign inverted. The threshold was selected by sweep-  
 629 ing [0.10, 0.50] on the calibration set and choosing the value that minimized misclassification of  
 630 non-sarcastic positive content.

### 631 C.4 Engagement weighting

$$w_e(t) = \min(\epsilon_0 + \log_{10}(\max(E(t), 1)), \epsilon_{\max}) \quad (9)$$

632 with  $\epsilon_0=0.5$ ,  $\epsilon_{\max}=3.0$ , and platform-specific engagement aggregates: Reddit/HN: `score+comments`;  
 633 Bluesky: `likes+reposts+replies`; Stack Overflow: `score+answers`;  
 634 Mastodon: `favourites+reblogs+replies`; GitHub Discussions: `reactions+replies`. The  
 635 logarithmic scaling prevents a single viral post from dominating an agent’s sentiment score.

## 636 D Sensitivity Analyses

637 This appendix documents the full sensitivity analyses including single-factor perturbations (sum-  
638 mary in main text), multi-factor perturbations, and bootstrap confidence intervals.

### 639 D.1 Single-factor perturbations

640 Table 5 in the main text shows rank changes for five representative agents under  $\pm 10$  percentage-  
641 point perturbations of each factor weight. The pattern across all 50 agents is consistent: no agent  
642 shifts by more than one rank position under any single-factor perturbation, and the SWE-category  
643 leader is invariant.

### 644 D.2 Multi-factor perturbations

645 We additionally test multi-factor perturbations: simultaneously varying two factor weights by  $\pm 5$ pp  
646 each (with the other two adjusted proportionally to preserve unit sum). Across  $\binom{4}{2}=6$  factor pairs  
647 and 4 sign combinations, we evaluate 24 multi-factor perturbations. The maximum rank shift ob-  
648 served across all 24 perturbations and all 50 agents is  $\pm 2$  ranks; no agent’s composite shifts by more  
649 than  $\pm 0.025$  in absolute score. The full perturbation matrix is included in the released artifact.

### 650 D.3 Bootstrap confidence intervals

651 We compute bootstrap confidence intervals on per-agent composite scores via 1,000 resamples of  
652 the underlying signal data. For each resample, we (i) draw with replacement from the collected  
653 text mentions per agent (preserving per-platform counts), (ii) recompute sentiment, (iii) recompute  
654 composite. Across 1,000 replicates, no agent in the top 20 shifts in median composite by more than  
655  $\pm 0.018$ . The 95% bootstrap intervals are tighter than the inter-agent score gaps for  $\geq 18$  of the top  
656 20 agents, supporting the robustness of the headline rankings.

### 657 D.4 Robustness to subset selection

658 A separate concern is whether the headline factor-independence and predictive-validity correlations  
659 are stable to which agents are included. We test this by leave-one-out resampling: for each of the  
660 34 public-repo agents, we drop that agent and recompute the inter-factor Spearman correlations and  
661 the headline IDE-bucket correlations. No single agent drives any of the three IDE-bucket results:  
662 B-only  $\rho_s$  ranges  $[0.42, 0.55]$  for VS Code installs (full-sample 0.48),  $[0.55, 0.65]$  for Open VSX  
663 (full-sample 0.60), and  $[0.46, 0.58]$  for JetBrains (full-sample 0.52). Maximum per-proxy deviation  
664 is  $|\Delta| \leq 0.07$ , and all three correlations remain statistically significant ( $p < 0.05$ ) under every drop.

## 665 E Cross-Factor Predictive Validity: Methodology Detail

666 This appendix documents the cross-factor predictive validity analysis (Section 5.2) in detail, includ-  
667 ing the choice of external proxies, the construction of the sub-composite, and robustness checks.

### 668 E.1 Why this analysis matters

669 A natural objection to any internally-aggregated composite score is that it has no external ground  
670 truth. AgentPulse aggregates 18 signals; rankings produced from that aggregation could in principle  
671 reflect nothing more than the aggregation rule itself. The cross-factor predictive validity analysis ad-  
672 dresses this directly: it asks whether a strict subset of the framework’s factors can predict signals the  
673 framework does not aggregate. If yes, the framework is capturing externally-observable structure; if  
674 no, the framework is internally consistent but externally vacuous.

### 675 E.2 Constructing the Benchmark+Sentiment sub-composite

676 The sub-composite is computed as

$$AP_{B+S}(a) = w'_B \cdot B(a) + w'_S \cdot S(a) \tag{10}$$

677 with  $w'_B = 0.35/(0.35 + 0.20) = 0.636$  and  $w'_S = 0.20/(0.35 + 0.20) = 0.364$  (the original  $w_B$   
678 and  $w_S$  renormalized to sum to 1). This preserves the relative weighting between Benchmark and  
679 Sentiment from the full composite. We exclude both Adoption and Ecosystem because they contain  
680 GitHub-derived signals (GitHub stars in Adoption; contributors and issue close rate in Ecosystem)  
681 that would mechanically correlate with any GitHub-derived target.

### 682 E.3 Choice of external proxies and pre-registration note

683 We test external proxies grouped into three theoretical buckets. The IDE-marketplace bucket was  
684 the primary pre-specified analysis. Within this bucket, only VS Code Marketplace installs was  
685 tested in the original draft submission; the Open VSX and JetBrains proxies were added post-hoc as  
686 replication on independent platforms after the initial submission. We mark these explicitly as post-  
687 hoc and apply Holm correction within the bucket; the library-reuse and community-engagement  
688 buckets are also post-hoc additions. We emphasize: across the 7+ proxies tested, the IDE-bucket  
689 result is the only one that shows positive, significant correlation; reporting only the positive findings  
690 would constitute selective reporting, which is why we report the full set including the discriminative-  
691 validity negatives and the community-engagement nulls.

#### 692 IDE marketplace bucket.

- 693 • *VS Code Marketplace installs (log-scaled)*. Number of installations of an agent’s VS Code  
694 extension. Available for 9 of the 34 public-repo agents with VS Code extensions; we  
695 substitute 0 for non-marketplace agents and compute over all 34.
- 696 • *Open VSX installs (log-scaled)*. The Eclipse-Foundation–hosted open-source VS Code  
697 extension registry, used by VSCodium and JetBrains Fleet. Available for 7 of 34; same  
698 zero-imputation.
- 699 • *JetBrains Marketplace downloads (log-scaled)*. Plugin downloads from the JetBrains plu-  
700 gin repository (IntelliJ, PyCharm, etc.). Available for 8 of 34; same zero-imputation.

#### 701 Library-reuse bucket.

- 702 • *GitHub Dependents*. Count of public repositories that depend on the agent’s package,  
703 scraped from `github.com/<repo>/network/dependents`. Available for 30 of 34  
704 public-repo agents; we report the raw subset (excluding 4 failed scrapes) rather than zero-  
705 imputing to avoid conflating "no public repo" with "no dependents."
- 706 • *PyPI downloads (log-scaled, last month)*. Available for 14 of 34 agents with published  
707 PyPI packages; we report the raw subset.

#### 708 Community-engagement bucket.

- 709 • *GitHub stars (log-scaled)*. Standard developer-attention metric. Available for 31 of 34.
- 710 • *Stack Overflow question volume*. Tagged-question counts. Available for all 34 (zero for  
711 untagged agents).
- 712 • *Discord member counts (log-scaled)*. Pulled via Discord’s invite-API endpoint. Available  
713 for 14 of 34 agents with publicly-listed Discord servers.

### 714 E.4 Secondary statistics and robustness checks

715 The IDE bucket replicates across all three platforms (Table 3 in main text). Discrete public-IDE-  
716 marketplace shipping-decision test (Mann–Whitney  $U$  on  $B$ -scores, present vs. absent): VS Code  
717  $U=42$ ,  $p=0.011$ ; Open VSX  $U=12$ ,  $p=0.0011$ ; JetBrains  $U=27$ ,  $p=0.004$ . All three remain signif-  
718 icant after Holm correction within the bucket. Spearman secondary statistics (B-only):  $\rho_s=+0.48$   
719 (VS Code),  $+0.60$  (Open VSX),  $+0.52$  (JetBrains). The library-reuse bucket shows the expected  
720 negative direction: B-only  $\rho_s=-0.30$  ( $p=0.11$ , directionally consistent but underpowered) and B+S  
721  $\rho_s=-0.38$  ( $p=0.04$ ) for GitHub Dependents on the  $n=30$  raw subset; PyPI downloads are weakly  
722 negative but not significant. The community-engagement bucket produces null results (all  $|\rho_s|<0.30$ ,  
723 all NS).

724 We performed four robustness checks on the IDE-bucket headline:

725 Robustness check 1: Mann–Whitney  $U$  confidence intervals (rank-biserial bootstrap, 1,000 resam-  
726 ples). VS Code:  $r_{rb}=+0.60$  [+0.27, +0.86]; Open VSX:  $r_{rb}=+0.86$  [+0.55, +1.00]; JetBrains:  
727  $r_{rb}=+0.71$  [+0.39, +0.93]. All three intervals exclude zero.

728 Robustness check 2: leave-one-out (Spearman B-only). VS Code:  $\rho_s$  ranges [0.42, 0.55], full-  
729 sample 0.48. Open VSX: [0.55, 0.65], full-sample 0.60. JetBrains: [0.46, 0.58], full-sample 0.52.  
730 Maximum per-proxy deviation  $|\Delta|\leq 0.07$ .

731 Robustness check 3: Pearson on log-scaled targets.  $r=0.45$  (VS Code),  $r=0.61$  (Open VSX),  
732  $r=0.49$  (JetBrains). The substantive conclusions are unchanged.

733 Robustness check 4: alternative predictors. B-only Spearman gives  $\rho_s=0.48, 0.60, 0.52$ ; B+S gives  
734 0.42, 0.55, 0.50; Sentiment alone is 0.01,  $-0.08, 0.12$  (near zero). Benchmark carries the predictive  
735 signal; Sentiment is non-harmful but not load-bearing in this snapshot.

## 736 E.5 Interpreting the discriminative-validity result

737 The negative correlation between B+S and GitHub Dependents ( $\rho_s=-0.38, p=0.04, n=30$ ; B-only  
738  $\rho_s=-0.30, p=0.11$ , directionally consistent but underpowered) is not a failure mode but a feature.  
739 Library-reuse metrics measure how often a project is imported as a dependency by other reposi-  
740 tories. Frameworks (LangGraph 38,792 dependents; Claude MCP 50,709; CrewAI 18,319; AutoGen  
741 4,107; Browser Use 2,510) dominate this metric because their entire purpose is to be reused; they  
742 typically lack published task benchmarks and therefore receive the  $B=0.5$  prior. Standalone agents  
743 that score highly on benchmarks (Claude Code  $B=0.82$ , Tabnine  $B=0.72$ , Continue  $B=0.74$ ) are  
744 not consumed as libraries by other repositories: developers *use* Claude Code to write code; they  
745 don't import `claude_code` into their own programs. The framework correctly distinguishes these  
746 two constructs. A composite that positively predicted library-reuse would be conflating “agent that  
747 completes tasks” with “library that is reused,” which would be a methodological flaw, not a strength.

## 748 E.6 What this analysis does not show

749 The analysis shows that the Benchmark factor predicts the discrete public-IDE-marketplace shipping  
750 decision across three independent platforms on the present snapshot, and is directionally negative for  
751 library-reuse and null for community-engagement metrics. It does not show: (a) that the framework  
752 predicts adoption across the full 50-agent registry including closed-source agents whose market-  
753 place presence is unobservable; (b) that the full composite (which adds Adoption and Ecosystem)  
754 predicts something further beyond Adoption-derived signals; (c) that the same correlations would  
755 hold on a future snapshot, since the agent ecosystem evolves rapidly; (d) that benchmark capability  
756 predicts install *volume* within the marketplace-present subset, the within-presence rank correlations  
757 are non-significant on  $n=6-8$ , and the released harness is designed to support testing this claim once  
758 the marketplace-present subset has matured to  $n\geq 25$ ; (e) that the Sentiment factor adds independent  
759 predictive power for IDE adoption in this snapshot; its retained value is at the aspect level (Ap-  
760 pendix C). A stronger validation would test whether the full composite predicts a target independent  
761 of all 18 signals, e.g., a developer-preference survey or commercial-deployment counts. We propose  
762 these as future work and welcome external replications.

## 763 F Reproducibility and Release

### 764 F.1 Release artifact structure

765 The artifact (anonymized URL provided in supplementary submission) is a single repository with  
766 the following structure:

- 767 • `collectors/`, 18 signal-collector implementations, one per signal source. Each collector  
768 subclasses a common `BaseCollector` interface with `collect()`, `rate_limit_config`,  
769 and `retry_policy`.
- 770 • `nlp/`, the NLP scoring stack: VADER wrapper, TextBlob wrapper, FinBERT wrapper,  
771 DistilBERT-SST2 wrapper, ensemble combiner, sarcasm detector, aspect lexicons.
- 772 • `registry/agents.json`, the 50-agent registry with provider attribution, primary work-  
773 load category, and signal-source assignments.

- 774 • `data/snapshots/`, hourly composite scores in Parquet format, one file per scoring run.
- 775 • `data/texts/`, all 15,000 NLP-scored texts in JSONL format, with composite quality  
776 scores, per-component sentiment, and aspect scores.
- 777 • `scripts/`, reproduction scripts that regenerate every table and figure from the released  
778 CSVs.
- 779 • `croissant.json`, Croissant metadata for machine-readable dataset documentation.
- 780 • `README.md`, setup, deployment, and reproduction instructions.
- 781 • `LICENSE`, CC BY 4.0 for the data and metadata; the code is released under MIT for com-  
782 patibility with downstream redistribution.

## 783 **F.2 Reproduction commands**

784 After cloning the repository and installing dependencies (`pip install -r requirements.txt`),  
785 `make all` regenerates every table and figure in the paper from the released CSVs (per-table scripts  
786 under `scripts/` are also enumerated in the repository `README`); no additional data collection is  
787 required.

## 788 **F.3 Compute requirements**

789 The full pipeline (collection, NLP, scoring) runs on a single commodity server with 8 CPU cores  
790 and 16 GB RAM. No GPU is required for any component (the neural sentiment models DistilBERT-  
791 SST2 and FinBERT run on CPU at acceptable throughput given the modest text volume). End-to-  
792 end reproduction of the paper’s tables and figures from the released CSVs takes under 5 minutes  
793 on consumer hardware. Live continuous deployment requires persistent network access for the  
794 collectors but no specialized hardware.

## 795 **F.4 License**

796 The released artifact is dual-licensed: the data and metadata under CC BY 4.0 (Creative Commons  
797 Attribution); the code under MIT. Both licenses permit redistribution and modification with attribu-  
798 tion. We chose CC BY for the data because it is the license recommended by NeurIPS for ED Track  
799 artifacts and is broadly compatible with downstream research use.

## 800 **F.5 Anonymization for review**

801 For the review period, the code and data are hosted on `anonymous.4open.science` and an  
802 anonymized data archive, respectively. URLs are provided in the supplementary submission. All  
803 commit history, author metadata, and embedded URLs have been scrubbed; reviewers can clone,  
804 run, and inspect the artifact without identifying the authors.

805 **NeurIPS Paper Checklist**

806 **1. Claims**

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

809 Answer: [Yes]

810 Justification: The abstract and Section 1 state three validation claims plus one diag-  
811 nostic, (i) factor-independence ( $n=50$ ,  $\rho_{\max}=0.75$  for Adoption–Ecosystem, all other  
812 pairwise  $|\rho|\leq 0.34$ ), (ii) circularity-controlled cross-factor predictive validity framed as  
813 the *public-IDE-marketplace shipping decision*, primary test Mann–Whitney  $U$  on Bench-  
814 mark for present vs. absent agents ( $n=34$ ; VS Code  $U=42$ ,  $p=0.011$ ,  $r_{rb}=+0.60$ ; Open  
815 VSX  $U=12$ ,  $p=0.0011$ ,  $r_{rb}=+0.86$ ; JetBrains  $U=27$ ,  $p=0.004$ ,  $r_{rb}=+0.71$ ; VS Code  
816 pre-specified, Open VSX and JetBrains post-hoc external replications), with Spearman  
817  $\rho_s=0.48$ – $0.60$  B-only ( $0.42$ – $0.55$  B+S) reported as a confirmatory secondary statistic, (iii)  
818 directional discriminative validity against library-reuse metrics (B-only  $\rho_s=-0.30$ , NS;  
819 B+S  $\rho_s=-0.38$ ,  $p=0.04$ ). The fourth analysis, ranking divergence between benchmark-  
820 only and composite ordering on the  $n=11$  SWE-bench overlap ( $\rho_s=0.09$ , 7 of 11 agents  
821 shifting  $\geq 2$  ranks), is presented as a diagnostic of the SWE-bench subset (Section 6)  
822 rather than as a positive validity claim. Claims (i)–(iii) are supported by Sections 5.1 and  
823 5.2; pre-registered falsifiers and the  $T+6$  replication commitment appear at the start of §5;  
824 remaining limitations are acknowledged in Section 7.

825 **2. Limitations**

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

827 Answer: [Yes]

828 Justification: Section 7 dedicates explicit paragraphs to: (a) the framework as a measure-  
829 ment tool not a single ground truth; (b) the closed-source measurement boundary; (c)  
830 sentiment selection bias and small-scale calibration; (d) underpowered  $n=11$  SWE-bench  
831 validation; (e) a falsifiability statement specifying conditions that would refute the frame-  
832 work. Appendix A documents three known failure modes of the data-quality pipeline. Ap-  
833 pendix C documents the limitation of small-scale internal annotation.

834 **3. Theory assumptions and proofs**

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

837 Answer: [N/A]

838 Justification: The paper makes no formal theoretical claims. All results are empirical:  
839 factor-independence correlations, cross-factor predictive validity, ranking divergence, fac-  
840 tor ablations, and weight-perturbation sensitivity.

841 **4. Experimental result reproducibility**

842 Question: Does the paper fully disclose all the information needed to reproduce the main  
843 experimental results to the extent that it affects the main claims and/or conclusions of the  
844 paper?

845 Answer: [Yes]

846 Justification: Section 3 provides full mathematical specifications of the four-factor compos-  
847 ite (Equations 1–5 with all weights and normalization constants); Section 4 documents the  
848 18-signal collection pipeline; Appendices A–F provide data-quality protocol parameters,  
849 per-platform statistics, the full 50-agent registry, NLP composite weights, aspect-lexicon  
850 construction, sensitivity analysis details, cross-factor predictive validity methodology, and  
851 reproduction scripts. The artifact is submitted in final form.

852 **5. Open access to data and code**

853 Question: Does the paper provide open access to the data and code, with sufficient instruc-  
854 tions to faithfully reproduce the main experimental results?

855 Answer: [Yes]

856 Justification: Appendix F documents the full open-source release under CC BY 4.0 (data)  
857 and MIT (code): collectors, NLP pipeline, agent registry, hourly score snapshots, NLP-  
858 scored texts, and per-table reproduction scripts. The release is accessible to reviewers via

859 an anonymized URL in the supplementary submission. The pipeline operates on free-tier  
860 public APIs for basic reproduction. A Croissant metadata file is included.

## 861 6. Experimental setting/details

862 Question: Does the paper specify all the training and test details necessary to understand  
863 the results?

864 Answer: [Yes]

865 Justification: All factor weights, normalization constants, decay parameters, and aggrega-  
866 tion functions are specified explicitly in Section 3. The data-quality threshold  $q^*=0.30$ ,  
867 sub-score weights, trigram-Jaccard threshold  $\tau=0.85$ , sentiment composite weights, and  
868 engagement-weighting bounds appear in Appendices A and C. Validation analyses report  
869 sample sizes ( $n=50$  for factor independence,  $n=34$  for IDE-bucket predictive validity  
870 with zero-imputation,  $n=30$  for library-reuse discriminative validity,  $n=11$  for SWE-bench  
871 overlap), test types (Mann–Whitney  $U$  primary for the IDE bucket; Spearman secondary  
872 throughout), Holm correction within the IDE bucket, and significance thresholds.

## 873 7. Experiment statistical significance

874 Question: Does the paper report error bars suitably and correctly defined or other appropri-  
875 ate information about the statistical significance of the experiments?

876 Answer: [Yes]

877 Justification: Table 3 reports the primary IDE-bucket test (Mann–Whitney  $U$  on Benchmark  
878 for marketplace-present vs. absent agents) with  $U$ ,  $p$ , and rank-biserial  $r_{rb}$  for each of  
879 the three platforms; Spearman  $\rho_s$  with 95% bootstrap CIs (1,000 resamples) for B-only  
880 and B+S is reported as a secondary statistic. Holm correction is applied within the IDE  
881 bucket. The factor-independence analysis ( $n=50$ ) is adequately powered. The SWE-bench-  
882 overlap analysis ( $n=11$ ) is explicitly framed as a diagnostic. Appendix D reports per-  
883 agent composite-score bootstrap intervals (median shift  $\leq \pm 0.018$  across the top 20), and  
884 Appendix E reports rank-biserial CI bootstrap, leave-one-out ( $\max |\Delta| \leq 0.07$ ), Pearson on  
885 log-scaled targets, and Dirichlet weight-perturbation bootstrap (1,000 draws) under which  
886 all three IDE-bucket correlations remain positive in every draw.

## 887 8. Experiments compute resources

888 Question: Does the paper provide sufficient information on the computer resources needed  
889 to reproduce the experiments?

890 Answer: [Yes]

891 Justification: Appendix F states the pipeline runs on a single commodity server (8 CPU  
892 cores, 16 GB RAM, no GPU); end-to-end reproduction of paper tables and figures from the  
893 released CSVs takes under 5 minutes on consumer hardware.

## 894 9. Code of ethics

895 Question: Does the research conducted in the paper conform, in every respect, with the  
896 NeurIPS Code of Ethics?

897 Answer: [Yes]

898 Justification: All data sources are public APIs and platform-terms-of-service-compliant  
899 collection (GitHub, PyPI, npm, VS Code Marketplace, Open VSX, JetBrains Marketplace,  
900 Discord invite API, Bluesky, Reddit, Hacker News, Stack Overflow, Mastodon, Dev.to,  
901 V2EX, Lemmy). No personally identifiable information is collected, retained, or released;  
902 sentiment scoring operates on aggregated text mentioning evaluated agents, not on individ-  
903 ual users. The framework evaluates AI agents, not human subjects.

## 904 10. Broader impacts

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

907 Answer: [Yes]

908 Justification: *Positive:* AgentPulse provides developers, researchers, and procurement  
909 teams a more complete signal-grounded view of agent quality than benchmark-only rank-  
910 ings, improving informed tool selection. *Negative:* as Section 7 notes, the framework struc-  
911 turally disadvantages closed-source agents whose adoption is not publicly observable; we  
912 make this measurement boundary explicit and recommend interpreting AgentPulse rank-  
913 ings involving such agents alongside benchmark-only rankings. Sentiment-based factors  
914 are subject to selection bias from over-represented vocal users.

- 915 **11. Safeguards**
- 916 Question: Does the paper describe safeguards that have been put in place for responsible  
917 release of data or models that have a high risk for misuse?
- 918 Answer: [N/A]
- 919 Justification: The released artifact is an evaluation framework and aggregated public-signal  
920 dataset, not a generative model, image dataset, or scraped personal-data corpus. Released  
921 sentiment scores are aggregated to the agent level and do not retain user identifiers. The  
922 artifact poses no identifiable misuse risk requiring safeguards.
- 923 **12. Licenses for existing assets**
- 924 Question: Are the creators or original owners of assets used in the paper properly credited  
925 and are the license and terms of use explicitly mentioned and properly respected?
- 926 Answer: [Yes]
- 927 Justification: VADER (5) and FinBERT (1) are cited and used under their respective open-  
928 source licenses (MIT, Apache 2.0); TextBlob (9) and the DistilBERT-SST2 sentiment  
929 model (13) are used as standard tools under their respective open-source licenses (MIT  
930 and Apache 2.0). All benchmark sources used by the composite (SWE-bench Verified (6),  
931 GAIA (10), WebArena (16), HumanEval+ (3), TAU-bench (14)) and the comparison bench-  
932 marks discussed in related work (AgentBench (8)) are cited and used per their published  
933 terms. Public-API data (GitHub, PyPI, npm, VS Code Marketplace, social platforms) is  
934 collected under each provider’s terms of service.
- 935 **13. New assets**
- 936 Question: Are new assets introduced in the paper well documented and is the documenta-  
937 tion provided alongside the assets?
- 938 Answer: [Yes]
- 939 Justification: The released AgentPulse artifact includes README documentation, schema  
940 files for all CSV/JSONL outputs, the 50-agent registry as JSON with category mappings, 18  
941 signal-collector specifications with cadences, the NLP scoring stack with domain-specific  
942 aspect lexicons, the four-factor composite implementation, per-table reproduction scripts,  
943 and a Croissant metadata file. Released under CC BY 4.0 (data) and MIT (code).
- 944 **14. Crowdsourcing and research with human subjects**
- 945 Question: For crowdsourcing experiments and research with human subjects, does the pa-  
946 per include the full text of instructions given to participants and screenshots, if applicable,  
947 as well as details about compensation?
- 948 Answer: [N/A]
- 949 Justification: AgentPulse evaluates AI agents using public observable signals. The only hu-  
950 man involvement was a 200-text annotation set for sentiment-pipeline calibration (Cohen’s  
951  $\kappa=0.81$ , two annotators drawn from the author team, no compensation as this was internal  
952 calibration). No external crowdsourcing or research with human subjects was conducted.  
953 Public social-media posts are collected under platform terms of service for sentiment ag-  
954 gregation, not as research-subject data.
- 955 **15. Institutional review board (IRB) approvals or equivalent for research with human  
956 subjects**
- 957 Question: Does the paper describe potential risks incurred by study participants, whether  
958 such risks were disclosed, and whether IRB approvals (or equivalent) were obtained?
- 959 Answer: [N/A]
- 960 Justification: The work does not involve research with human subjects requiring IRB re-  
961 view. Public-signal aggregation from platform APIs operating under their published terms  
962 of service does not constitute human-subjects research under standard institutional defini-  
963 tions; the internal 200-text calibration annotation involved authors only.
- 964 **16. Declaration of LLM usage**
- 965 Question: Does the paper describe the usage of LLMs if it is an important, original, or  
966 non-standard component of the core methods?
- 967 Answer: [N/A]

968 Justification: LLMs are not a methodological component of the AgentPulse framework.  
969 The framework evaluates LLM-based AI agents but does not itself use LLMs for scor-  
970 ing, classification, or aggregation; the NLP pipeline uses lightweight transformer models  
971 (FinBERT, DistilBERT-SST2) and lexicon-based methods (VADER, TextBlob), all cited as  
972 standard tools. Any LLM use was limited to writing assistance and does not affect method-  
973 ology, scientific rigor, or originality.