

# Early Warning of High-Cost Agentic Pull Requests via Scenario–Cost Modeling

Anonymous Author(s)

## Abstract

Using the AIDev dataset (*snapshot/version: [FILL IN]*), we analyze 33,596 agentic PRs and (RQ1) categorize them into interaction scenarios: S0 (solo agent) 32.94%, S1 (human reviewed) 12.52%, and S2 (human co-edited) 54.55%. We then (RQ2) define a composite cost model spanning review intensity, communication, and iteration, and (RQ3) predict high-cost PRs under a fixed alert-budget policy.

## CCS Concepts

- Software and its engineering → Software maintenance tools.

## Keywords

Mining Software Repositories, agentic pull requests, code review, cost modeling, early warning

## ACM Reference Format:

Anonymous Author(s). 2025. Early Warning of High-Cost Agentic Pull Requests via Scenario–Cost Modeling. In *Proceedings of MSR ’26: Proceedings of the 23rd International Conference on Mining Software Repositories (MSR 2026)*. ACM, New York, NY, USA, 3 pages. <https://doi.org/10.1145/nmnnnnnn>.

## 1 Introduction

Agentic PRs can reduce developer effort, yet maintainers may face increased review burden and coordination overhead. This paper studies whether simple, early-available signals can provide actionable warnings about high-cost PRs, enabling maintainers to triage limited review resources.

## 2 Dataset and Experimental Setup

*Dataset.* We use the AIDev dataset (*snapshot/version: [FILL IN] EXACT SNAPSHOT, DATE, OR COMMIT*). Because the dataset is continuously updated, we report the exact snapshot used for all analyses.

*Unit of analysis.* Our unit is the PR. We analyze PR metadata, review events, and comments as available in AIDev tables (e.g., `pull_request`, `pr_reviews`, `pr_comments`, `pr_review_comments_v2`).

*Reproducibility.* We will release a replication package including SQL extraction scripts, analysis notebooks, and figure-generation code (*link/DOI: [FILL IN]*).

---

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from [permissions@acm.org](mailto:permissions@acm.org).

MSR 2026, Rio de Janeiro, Brazil

© 2025 Copyright held by the owner/author(s). Publication rights licensed to ACM. ACM ISBN 978-x-xxxx-xxxx-x/YY/MM  
<https://doi.org/10.1145/nmnnnnnn.nmnnnnnn>

**Table 1: Scenario distribution by agent (percent within each agent).**

| Agent        | S0     | S1     | S2     |
|--------------|--------|--------|--------|
| Claude_Code  | 33.12% | 11.76% | 55.12% |
| Copilot      | 28.71% | 41.37% | 29.92% |
| Cursor       | 50.88% | 18.23% | 30.89% |
| Devin        | 39.65% | 31.92% | 28.42% |
| OpenAI_Codex | 31.14% | 1.25%  | 67.61% |

**Table 2: High-cost rate by interaction scenario.**

| Scenario             | n     | High-cost n | Rate   |
|----------------------|-------|-------------|--------|
| S0 (Solo agent)      | 11065 | 549         | 0.0496 |
| S1 (Human reviewed)  | 4205  | 2551        | 0.6067 |
| S2 (Human co-edited) | 18326 | 3619        | 0.1975 |

## 3 RQ1: Interaction Scenarios of Agentic PRs

*Goal.* Characterize how agentic PRs differ in human involvement.

*Scenario definitions.* We label each PR into one of three scenarios: (i) **S0 (Solo agent)**: no human comments/reviews/feedback; merging by humans is allowed. (ii) **S1 (Human reviewed)**: human comments/reviews exist, but no human commits. (iii) **S2 (Human co-edited)**: at least one human-authored commit exists.

*Results.* Across all agents (N=33,596 PRs), S0 accounts for 11,065 PRs (32.94%), S1 for 4,205 PRs (12.52%), and S2 for 18,326 PRs (54.55%). Scenario distributions vary substantially across agents. For example, Copilot has the highest share of S1 (41.37%), while OpenAI\_Codex is dominated by S2 (67.61%) with a very small S1 share (1.25%).

*Outputs.* We produce a scenario-labeled PR table `pr_scenarios_rq1` used in later RQs.

## 4 RQ2: Cost Model for Review, Communication, and Iteration

*Goal.* Quantify PR cost along multiple dimensions relevant to maintainers.

*Cost components.* We define three cost dimensions: (1) **Review intensity**: `review_count` and `request_changes_count`; (2) **Communication cost**: `comment_count` (including review comments); (3) **Iteration cost**: `post_review_review_count` (number of review rounds after the first review; used as an iteration proxy).

*High-cost label.* We compute a cost score via log-summed components and label *high\_cost* as the top 20% of PRs by cost score. In our dataset, 6,719 out of 33,596 PRs are labeled as high-cost (20.00%).

117 **Table 3: Early-warning performance under fixed alert bud-**  
 118 **gets (repo-level split).**

| Alert budget | $k$  | AUC   | Precision@ $k$ | Recall@ $k$ |
|--------------|------|-------|----------------|-------------|
| Top-10%      | 1301 | 0.831 | 0.694          | 0.727       |
| Top-20%      | 2601 | 0.831 | 0.409          | 0.855       |

124 **Table 4: High-cost rate by agent and scenario (cell shows rate**  
 125 **with sample size).**

| Agent        | S0 (Solo)         | S1 (Human reviewed) | S2 (Human co-edited) |
|--------------|-------------------|---------------------|----------------------|
| Claude_Code  | 0.06 ( $n=152$ )  | 0.33 ( $n=54$ )     | 0.49 ( $n=253$ )     |
| Copilot      | 0.01 ( $n=1427$ ) | 0.64 ( $n=2056$ )   | 0.81 ( $n=1487$ )    |
| Cursor       | 0.16 ( $n=784$ )  | 0.50 ( $n=281$ )    | 0.61 ( $n=476$ )     |
| Devin        | 0.13 ( $n=1914$ ) | 0.65 ( $n=1541$ )   | 0.65 ( $n=1372$ )    |
| OpenAI_Codex | 0.02 ( $n=6788$ ) | 0.26 ( $n=273$ )    | 0.08 ( $n=14738$ )   |

## 5 RQ3: Early Warning of High-Cost PRs

136 *Goal.* Predict whether an incoming PR will be high-cost using  
 137 only early-available signals.

139 *Features and split.* We train a logistic regression classifier with  
 140 categorical early signals: *agent*, *scenario label* (RQ1), and *PR state*.  
 141 We evaluate using a **repository-level split** to reduce within-repo  
 142 leakage and report AUC.

144 *Budget-based alerting (Top- $k$ ).* Because the feature space is coarse-  
 145 grained and yields tied risk scores, probability thresholding can  
 146 produce unstable alert volumes. We instead use a fixed alert-budget  
 147 policy: flag only the Top- $k$  highest-risk PRs in each batch (ties broken  
 148 deterministically by PR id), matching practical review-resource  
 149 constraints.

150 *Results.* The model achieves AUC = 0.831. Under a Top-10%  
 151 alert budget ( $k = 1301$ ), it attains Precision@10% = 0.694 and  
 152 Recall@10% = 0.727 (F1 = 0.710). Under a Top-20% budget ( $k =$   
 153 2601), it achieves Precision@20% = 0.409 and Recall@20% = 0.855  
 154 (F1 = 0.553), demonstrating a clear precision–coverage trade-off  
 155 (Table 3).

157 *Interpretability.* Agent-level high-cost rates remain well sepa-  
 158 rated after accounting for uncertainty. For example, Copilot has a  
 159 high-cost rate of 0.5119 (Wilson 95% CI [0.4980, 0.5258]) and Devin  
 160 0.4411 ([0.4271, 0.4551]), substantially higher than OpenAI\_Codex  
 161 0.0617 ([0.0585, 0.0649]), supporting agent identity and scenario  
 162 labels as actionable early-warning features.

164 *Agent × scenario heterogeneity.* High-cost risk is highly hetero-  
 165 geneous across the intersection of agent and scenario. For example,  
 166 Copilot exhibits extremely high high-cost rates in S2 (0.81,  $n=1487$ )  
 167 and S1 (0.64,  $n=2056$ ), but remains very low in S0 (0.01,  $n=1427$ ).  
 168 In contrast, OpenAI\_Codex shows a moderate high-cost rate in  
 169 S1 (0.26,  $n=273$ ) but a much lower rate in S2 (0.08,  $n=14738$ ), sug-  
 170 gesting that the “human-reviewed” workflow is disproportionately  
 171 associated with high-cost PRs for certain agents. These patterns sup-  
 172 port scenario label and agent identity as actionable early-warning  
 173 signals.



197 **Figure 1: High-cost rate by agent with Wilson 95% confidence**  
 198 **intervals.**

## 6 Discussion and Implications

201 Our findings suggest maintainers can control alert noise versus  
 202 coverage by choosing an alert budget. A Top-10% policy yields  
 203 high precision (fewer false alarms), while Top-20% captures most  
 204 high-cost PRs. The strong concentration of high-cost PRs in human-  
 205 reviewed workflows indicates that collaboration mode and agent  
 206 identity are meaningful early signals for triage.

## 7 Ethical Implications

209 We analyze publicly available repository artifacts and report aggre-  
 210 gate results. We avoid releasing any personally identifying infor-  
 211 mation (PII) and do not attempt to deanonymize users. Automated  
 212 warnings may influence maintainers’ attention; therefore, alerts  
 213 should be used as decision support rather than as automated rejec-  
 214 tion signals, and should be periodically audited for unintended bias  
 215 across projects or contributors.

## 8 Threats to Validity

219 *Construct validity.* Our cost model uses observable proxies (re-  
 220 view rounds, comments, request-changes) and may not capture all  
 221 forms of effort (e.g., offline discussion).

225 *Internal validity.* Some event logs may be incomplete; we mit-  
 226 igate this by relying on stable tables for review/comment counts  
 227 and by using repository-level splits.

228 *External validity.* Results are specific to the AIDev snapshot and  
 229 the studied repositories; agent distributions are imbalanced, and  
 230 small-sample cells should be interpreted cautiously.

## 233 9 Conclusion

234 We propose a scenario–cost framework for agentic PRs and show  
 235 that early-available categorical signals can provide strong high-cost  
 236 risk ranking ( $AUC = 0.831$ ) and practical Top- $k$  early warning under  
 237 fixed alert budgets. Future work will incorporate richer early PR  
 238

239 features (e.g., code-change characteristics) and study downstream  
 240 outcomes such as acceptance and turnaround time.

## 241 References

|     |     |
|-----|-----|
| 242 | 291 |
| 243 | 292 |
| 244 | 293 |
| 245 | 294 |
| 246 | 295 |
| 247 | 296 |
| 248 | 297 |
| 249 | 298 |
| 250 | 299 |
| 251 | 300 |
| 252 | 301 |
| 253 | 302 |
| 254 | 303 |
| 255 | 304 |
| 256 | 305 |
| 257 | 306 |
| 258 | 307 |
| 259 | 308 |
| 260 | 309 |
| 261 | 310 |
| 262 | 311 |
| 263 | 312 |
| 264 | 313 |
| 265 | 314 |
| 266 | 315 |
| 267 | 316 |
| 268 | 317 |
| 269 | 318 |
| 270 | 319 |
| 271 | 320 |
| 272 | 321 |
| 273 | 322 |
| 274 | 323 |
| 275 | 324 |
| 276 | 325 |
| 277 | 326 |
| 278 | 327 |
| 279 | 328 |
| 280 | 329 |
| 281 | 330 |
| 282 | 331 |
| 283 | 332 |
| 284 | 333 |
| 285 | 334 |
| 286 | 335 |
| 287 | 336 |
| 288 | 337 |
| 289 | 338 |
| 290 | 339 |
| 291 | 340 |
| 292 | 341 |
| 293 | 342 |
| 294 | 343 |
| 295 | 344 |
| 296 | 345 |
| 297 | 346 |
| 298 | 347 |
| 299 | 348 |