

# DATA 542 Project – Milestone 2: Mid-Project Report

HAOZHONG JI, University of British Columbia Okanagan, Canada

YIN-WEN TSAI, University of British Columbia Okanagan, Canada

## 1 Research Questions and Dataset

We use the AIDev dataset. It is about pull requests on GitHub that use AI coding agents. The dataset has information on pull requests, repositories, users, reviews, and commits. In this project, we only study pull requests that involve AI agents. We also keep only pull requests with valid times. Their final status must be clear. So they are either merged or closed.

We will answer three research questions (RQs). Each RQ uses more than one feature.

- **RQ1:** Do different AI coding agents get used in different types of repositories? We compare repositories with different popularity and activity.
- **RQ2:** How do the size and type of AI pull requests relate to review work and merge results, when we group pull requests into simple size and file-type buckets?
- **RQ3:** For AI-agent pull requests, how do simple human–AI collaboration patterns (based on follow-up commits) relate to merge rate and time-to-merge?

## 2 Methodology

For all RQs, we will load the pull request table and the repository table. We will join them by repository ID. We will keep only AI-related pull requests with clear outcomes. We will create features for our analysis. These include agent type, repository stars, and repository activity. We will also use commit-level tables and user information to detect human and agent commits inside each pull request.

**RQ1 (adoption patterns).** We will group repositories by star counts. For example, low, medium, and high stars. We will count how many AI pull requests each agent makes in each group. We will also look at changes over time, such as by month. We will show the results with line charts and bar charts.

**RQ2 (characteristics and outcomes).** For each pull request, we will measure its size. We use lines added, lines deleted, and files changed. We will label file types using file paths. For example, code, tests, docs, or config files. We will record review comments, merge outcome, and time-to-merge. We will then group pull requests into simple size buckets (for example, small, medium, large) and by main file type. For each group, we will summarize merge rate, median time-to-merge, and average number of review comments. We will compare these groups using tables, bar charts, and box plots. If time allows, we may run a very simple logistic regression as an exploratory check, but our main focus will be on data wrangling and descriptive analysis.

**RQ3 (collaboration patterns).** For RQ3, we will use commit-level tables joined with the user table to detect human and agent commits inside each pull request. In AIDev, every pull request is opened by an AI agent, so we study teamwork only through follow-up commits. For each pull request, we will count human and agent commits after the first agent commit and assign a simple pattern label, such as “agent-only” (no human commits) or “agent-then-human” (at least one human follow-up commit). For each pattern, we will summarize merge rate, median time-to-merge, and average review comments and compare the patterns with tables and a few bar or box plots.

---

Authors’ Contact Information: Haozhong Ji, jhzs@student.ubc.ca, University of British Columbia Okanagan, Kelowna, Canada; Yin-Wen Tsai, ywtsai@student.ubc.ca, University of British Columbia Okanagan, Kelowna, Canada.

### 3 Results

#### 3.1 Results for RQ1

**3.1.1 Data and grouping.** For RQ1, we use the AIDev dataset and focus on pull requests (PRs) opened by AI coding agents. Each PR is linked to its repository, so we can attach repository-level features such as stars and activity. We only keep PRs with a clear final status (merged or closed) and valid timestamps.

To measure *popularity*, we use the number of stars on each repository. We rank repositories by star count and then cut the rank into three equal-sized groups: *Low*, *Medium*, and *High* popularity. This rank-based grouping helps reduce the effect of a few extremely famous projects.

To measure *activity*, we build a simple activity score from repository-level events (for example, commits, issues, and pull requests). We then split this score into three groups: *Low*, *Medium*, and *High* activity.

For both popularity and activity, we count how many AI-generated PRs each agent creates in each group. Then we convert these counts to proportions within each group. This gives two main tables:

- popularity group  $\times$  agent (pop\_agent),
- activity group  $\times$  agent (act\_agent).

We also visualize these distributions using grouped bar charts, with one panel for popularity and another for activity.

All percentages in the report are rounded to one decimal place. This means that very small but non-zero shares (for example, less than 0.05%) can appear as 0.0% in the tables and plots. In those cases, we describe them as “negligible but non-zero”, since they exist in the raw counts but do not change the main patterns.

**3.1.2 Popularity and agent usage.** Overall, **OpenAI\_Codex** dominates AI usage in all popularity groups, but the mix of agents changes with popularity.

In *low-popularity* repositories, OpenAI\_Codex accounts for about 90.8% of AI PRs (260,202 PRs). Copilot contributes about 8.4% (24,193 PRs), and Claude\_Code about 0.8% (2,247 PRs). Other agents such as Cursor and Devin only appear, if at all, at negligible levels in this group. Their shares are so small (below 0.1%) that they round to 0.0% in our summary tables, so they do not affect the overall picture.

In *medium-popularity* repositories, the distribution becomes even more concentrated. Here, essentially all AI PRs (around 100%) are opened by **OpenAI\_Codex** (286,642 PRs). Any other agents in this group, if present, contribute at most a very small number of PRs. Their shares are too small to be visible at one-decimal rounding and can be treated as negligible for this analysis.

In *high-popularity* repositories, the picture becomes more diverse. OpenAI\_Codex still leads, but its share drops to about 74.7% (214,188 PRs). Copilot contributes about 5.6% (16,022 PRs), and Claude\_Code about 0.9% (2,487 PRs). This is also where Cursor and Devin reach noticeable levels: about 9.0% (25,880 PRs) for Cursor and 9.8% (28,066 PRs) for Devin. There may be a few additional agents with very small shares (again below 0.1%), but they are negligible compared with these main tools.

These results suggest that new or alternative AI agents (such as Cursor and Devin) are mainly used in more popular projects. Less popular repositories rely almost entirely on OpenAI\_Codex, with some additional usage of Copilot and Claude\_Code, and only negligible use of other agents.

**3.1.3 Activity and agent usage.** The patterns by *activity* are similar, but not identical.

In *low-activity* repositories, OpenAI\_Codex again dominates with about 88.7% of AI PRs (254,194 PRs). Copilot has about 10.3% (29,414 PRs), and Claude\_Code around 1.1% (3,034 PRs). Other agents,

including Cursor and Devin, only appear (if at all) at extremely small levels. Their shares are below our rounding threshold, so they show as 0.0% in the tables.

In *medium-activity* repositories, the distribution is again highly concentrated. OpenAI\_Codex accounts for essentially all AI PRs (about 100%, 286,642 PRs). Any other agents in this group have either zero or a very small number of PRs. Their shares are negligible compared with OpenAI\_Codex and do not change the main conclusion.

In *high-activity* repositories, usage becomes more mixed. OpenAI\_Codex still has the largest share, at about 76.8% (220,196 PRs). Copilot accounts for about 3.8% (10,801 PRs), and Claude\_Code about 0.6% (1,700 PRs). As with popularity, Cursor and Devin appear only in this high-activity group at substantial levels, with shares around 9.0% (25,880 PRs) and 9.8% (28,066 PRs), respectively. Again, any additional agents beyond these have very small shares and can be ignored at the scale of our plots.

So, repositories with higher activity levels are more likely to use a wider range of AI agents. Low- and medium-activity projects mostly depend on OpenAI\_Codex, with some Copilot and Claude\_Code and only negligible use of other tools.

**3.1.4 Summary of RQ1 Findings.** Across both popularity and activity, **OpenAI\_Codex** is the main AI agent in the AIDev dataset. It accounts for almost all AI-generated PRs in medium-popularity and medium-activity repositories, and it is still the majority agent in the other groups.

However, we observe clear differences at the high end. Highly popular and highly active repositories use a more diverse set of agents. In those groups, Cursor and Devin appear with non-trivial shares (around 9–10%), and the relative share of OpenAI\_Codex decreases, even though it remains dominant. Some other agents also appear at very small proportions, but their contributions are below 0.1% and do not change the main trends.

These patterns suggest that large or busy projects are more willing, or more able, to experiment with multiple AI tools. Smaller or less active repositories mostly adopt a single, established agent. This answers RQ1: different AI coding agents are indeed used in different types of repositories, and the diversity of AI tools increases in more popular and more active projects.

## 3.2 Results for RQ2

To investigate how PR size and dominant file type relate to merge outcomes and review effort, we analyzed merge rate, median time-to-merge, and the average number of review comments across three PR size buckets (Small, Medium, Large) and five dominant file types (code, config, docs, other, test). The results reveal consistent patterns across all metrics.

**3.2.1 Merge Rate.** According to Figure 1, small PRs show the highest merge rates across all file types, typically ranging from 85% to 100%. Test-dominated PRs exhibit the strongest merge likelihood, while code- and other-related PRs merge slightly less frequently but still maintain high overall acceptance.

Merge rates decline for Medium PRs and drop further for Large PRs. Large PRs dominated by test or other files show the lowest merge rates (around 60–73%), suggesting that reviewers are more cautious with large, complex contributions.

Our key observation is that: small PRs are consistently more mergeable. Merge rates decrease as PR size grows, regardless of file type.

**3.2.2 Time-to-Merge.** Small PRs in Figure 2 merge rapidly, often within minutes or hours. Medium PRs require slightly more review time, particularly for config-heavy changes.

Large PRs demonstrate the most substantial delays. In particular, PRs dominated by configuration files exhibit the longest median time-to-merge ( $\approx 0.19$  days). Large code PRs also face prolonged review, reflecting the higher verification burden associated with these changes.

We get the key observation: merge time increases with PR size, with config and code changes demanding the longest review effort for Large PRs.

**3.2.3 Review Comments.** Review activity follows a similar size-dependent trend from Figure 3. Small PRs receive roughly 2.7–2.9 comments on average, except test PRs, which receive significantly fewer (0.7 comments), indicating lower review friction.

Medium PRs receive more comments overall, with documentation PRs receiving the most discussion. For Large PRs, review intensity increases sharply. Config, code, and documentation PRs all receive between 3.8–4.3 comments on average, suggesting that reviewers allocate more effort to evaluating system-critical or widely visible changes.

The key observation of this part is: review effort scales with PR size, and config/code/docs PRs attract the most reviewer scrutiny.

**3.2.4 Summary of RQ2 Findings.** Across merge rate, merge time, and review comments, results align around two consistent insights:

- (1) PR size is the strongest determinant of merge behavior. Smaller PRs merge more successfully, more quickly, and with fewer comments.
- (2) Dominant file type modifies reviewer expectations and workload. Test and documentation changes are merged easily, while configuration and code changes—especially large ones—trigger slower merges and more review activity.

These findings suggest that contributing smaller, well-scoped PRs and isolating risky file types (such as config or core code) can substantially improve merge outcomes in AI-related repositories.

## 4 Interpretation

### 4.1 RQ1 Interpretation

### 4.2 RQ2 Interpretation

The results indicate that both PR size and dominant file type meaningfully influence merge behavior. Across all file types, merge rates consistently decrease as PR size increases, suggesting that reviewers are more hesitant to approve large AI-generated changes. This pattern aligns with general software engineering practices, where larger PRs carry higher review cost and perceived risk.

File type further differentiates review outcomes. Configuration-dominant PRs show the slowest time-to-merge and the highest review activity, likely because configuration errors can have system-wide impact and therefore require more careful validation. In contrast, documentation PRs are merged more quickly and with fewer comments, consistent with their lower risk profile. Small test PRs also show very high merge rates, indicating low resistance to accepting additional tests.

Overall, AI-generated PRs do not bypass typical reviewer behaviors. Instead, they reinforce existing patterns:

- (1) small, focused changes are easier to merge,
- (2) certain file types (e.g., configuration) trigger more scrutiny,
- (3) review workload scales with PR size.

These observations suggest that the effectiveness of AI-generated contributions depends not only on content quality but also on how the changes are structured and scoped.

5    **Link to GitHub Repository**  
data 542 project

6    **Remaining Work to Be Completed**

We will add the interpretation section for RQ1, complete the full analysis and reporting for RQ3, and finalize the code repository. The final project report will also be refined to ensure cohesion and consistency across all research questions.

A    **Appendix**

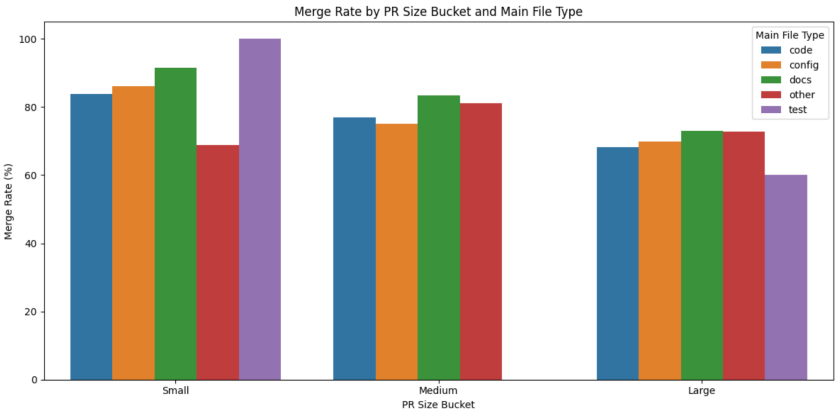


Fig. 1. Merge Rate by PR Size Bucket and Main File Type

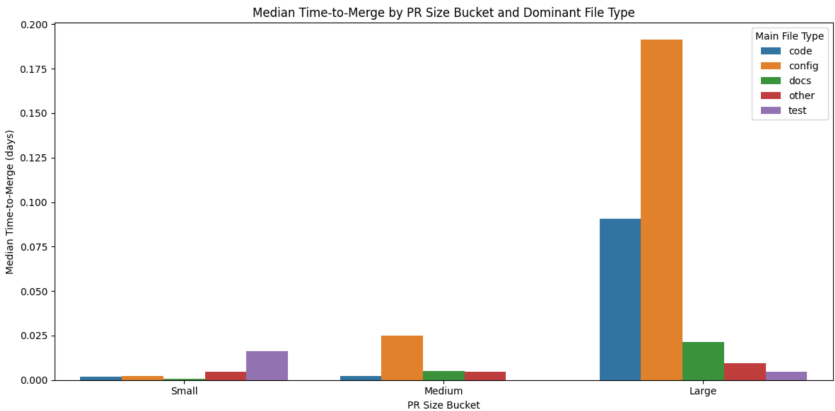


Fig. 2. Median Time-to-Merge by PR Size Bucket and Dominant File Type

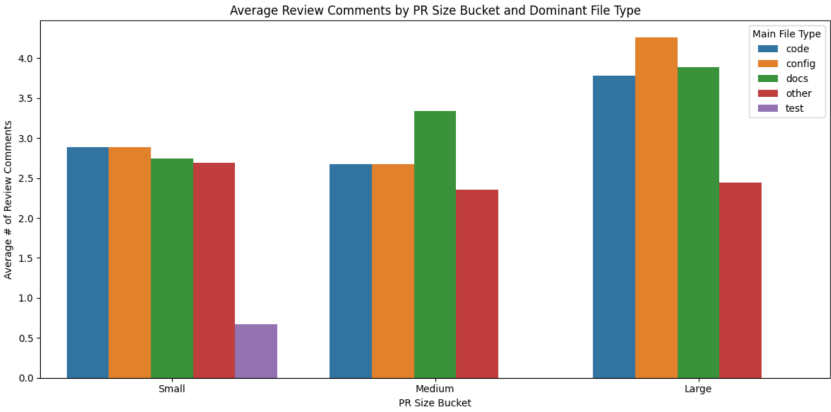


Fig. 3. Average Review Comments by PR Size Bucket and Dominant File Type