

Project Milestone 2 Report

YAOJUN LIU, University of British Columbia, Canada
SASIVIMOL SIRIJANGKAPATTANA, University of British Columbia, Canada

1. Introduction

This Milestone 2 report documents progress on the first research question (RQ1) of the project. The objective is to examine whether repository- and user-level characteristics predict pull request (PR) resolution time.

The project investigates three research questions:

- **RQ1:** How do repository popularity (stars, forks) and user influence (followers, account age) relate to PR resolution time?
- **RQ2:** How do text characteristics of PRs (length and sentiment) vary across repositories of different popularity tiers and users of different influence levels?
- **RQ3:** How does AI coding agent contribution behavior vary across repositories with different popularity levels?

This milestone focuses exclusively on RQ1. RQ2 and RQ3 will be completed for the final report.

2. Data Wrangling and Feature Construction

Three datasets were used: pull requests, repositories, and users. All timestamps were converted to datetime format and merged using `repo_id` and `user_id`. Only closed PRs were retained to compute resolution time.

2.1 Resolution Time

PR resolution time was computed as:

$$resolution_time_hours = \frac{closed_at - created_at}{3600}$$

User account age was defined as:

$$user_account_age_days = created_at(PR) - created_at_user.$$

2.2 Zero-Inflated Predictors

Repository popularity (stars, forks) and user influence (followers) are heavily zero-inflated (60–80% zeros). To preserve dataset representativeness, zero values were retained. All predictors and the outcome were log-transformed using $\log(1 + x)$ to reduce skew.

3. Preliminary Results for RQ1

3.1 Correlation Analysis

A correlation matrix of log-transformed variables shows weak associations between predictor variables and PR resolution time:

- $\text{corr}(\log(\text{resolution time}), \log(\text{stars})) \approx 0.20$
- $\text{corr}(\log(\text{resolution time}), \log(\text{forks})) \approx 0.22$
- $\text{corr}(\log(\text{resolution time}), \log(\text{followers})) \approx 0.079$

Authors' Contact Information: Yaojun Liu, University of British Columbia, Department of Computer Science, Okanagan, BC, Canada; Sasivimol Sirijangkappattana, University of British Columbia, Department of Computer Science, Okanagan, BC, Canada.

- 50 • $\text{corr}(\log(\text{resolution time}), \log(\text{user age})) \approx 0.038$

51 All correlations are below 0.25, indicating negligible linear relationships between popularity/in-
52 fluence and PR processing time.

54 3.2 Regression Analysis

55 A log-linear regression model was fitted using 817,668 closed pull requests:

57
58 $\log(1+Y) = \beta_0 + \beta_1 \log(1+\text{stars}) + \beta_2 \log(1+\text{forks}) + \beta_3 \log(1+\text{followers}) + \beta_4 \log(1+\text{user age}) + \varepsilon.$

59 **Model performance:**

61
62 $R^2 = 0.049, \quad \text{Adjusted } R^2 = 0.049.$

63 Thus, only about 4.9% of the variability in PR resolution time is explained by these predictors.

64 **Estimated coefficients (all $p < 0.001$):**

- 65 • β_1 (log-stars): +0.0163
- 66 • β_2 (log-forks): +0.2223
- 67 • β_3 (log-followers): -0.0120
- 68 • β_4 (log-account age): +0.0099

70 Although statistically significant due to the large sample size, these effects are extremely small
71 in magnitude.

72 3.3 Interpretation

73 These findings indicate that:

74
75 • Repository popularity has little to no predictive power for PR resolution time.

76 • User influence (followers, account age) also has negligible predictive impact.

77 • The small coefficient sizes and low R^2 suggest that PR resolution time is primarily driven
78 by unobserved factors such as maintainer availability, CI outcomes, review workload, and
79 project-specific processes.

80 This aligns with prior software engineering studies reporting that PR review speed is difficult to
81 predict from repository- or user-level metadata alone.

83 4. Limitations

84
85 • Resolution time is highly skewed, with many short-duration PRs and a long tail of very
86 long ones.

87 • Important complexity indicators (changed files, additions, deletions, comments) were not
88 used in this milestone but may explain additional variance.

89 • Zero-inflation limits the predictive strength of popularity and influence metrics.

90 5. Next Steps: RQ2 and RQ3

91 RQ2: How do text characteristics of pull requests vary across repositories of different 92 popularity levels and users of different influence levels?

93
94 We will analyze whether PR text features—such as sentiment and text length—differ across reposi-
95 tories with varying popularity (low, medium, high) and across users with different influence levels
96 (based on follower count). This analysis will help determine whether more popular repositories or
97 high-influence users tend to produce PRs with distinct linguistic or stylistic characteristics.

99 *Methodology.* PR, repository, and user data are merged using `repo_id` and `user_id`. Text-based
100 features are computed by concatenating PR titles and bodies into a single field, from which
101 `text_length` and sentiment polarity are derived using the `TextBlob` library. Repository popularity
102 tiers are generated using star count quantiles, while user influence tiers are based on follower
103 quantiles. Exploratory visualizations (e.g., boxplots) will then be used to compare sentiment and
104 text length across popularity and influence levels.

105 *Features Used.* The analysis will rely on the following fields:

- 106 • `title`, `body` (raw PR text)
- 107 • `full_text` (combined text used for NLP features)
- 108 • `text_length` (character length of PR text)
- 109 • `sentiment` (sentiment polarity extracted via `TextBlob`)
- 110 • `stars` (repository popularity)
- 111 • `followers` (user influence)
- 112 • `repo_id`, `user_id` (merge keys)

114 **RQ3: How do AI agents behave across repositories of different popularity levels?**

115 We will analyze how AI-generated pull requests differ when submitted to repositories with varying
116 popularity levels. The focus is on understanding whether AI agents contribute differently in terms
117 of submission volume, merge outcomes, and PR lifetime across low, medium, and high-popularity
118 repositories.

119 *Methodology.* PR metadata will be merged with repository information using `repo_id`. Repository
120 popularity will be categorized into three tiers (low, medium, high) based on star count quantiles. For
121 all PRs authored by AI agents, we will compute behavioral metrics—AI PR volume, closure/merge
122 rate, and PR lifetime (`closed_at` – `created_at`)—within each popularity tier. Exploratory data analysis,
123 including summary statistics and visualizations, will be used to compare AI behavior across these
124 tiers.

125 *Features Used.* The analysis will rely on the following fields:

- 126 • `agent` (AI-agent: Claude_Code, Copilot,..)
- 127 • `repo_id` (repository identifier)
- 128 • `stars` (repository popularity measure)
- 129 • `state` (merged/closed/open)
- 130 • `created_at` and `closed_at` (timestamps used to compute PR lifetime)
- 131 • `id` (pull request identifier)

132 **6. GitHub Repository**

133 Code and data processing pipeline:

134 <https://github.com/liuyaojun1/542.git>

135

136

137

138

139

140

141

142

143

144

145

146

147