

1 Factors Driving Developer Influence and Efficiency in the AI 2 Community

3 ZIHUAN LUI, University of British Columbia, Canada

4 AVALVIR KAUR SEKHON, University of British Columbia, Canada

5 We analyze the AIDev dataset to understand the factors driving developer influence and efficiency in the AI
6 development community. Focusing on user activity, repository metadata, and issue tracking, we investigate the
7 relationship between contribution volume and follower count, predictors of popularity, and factors affecting
8 issue resolution time. Our preliminary results indicate that contribution volume has a negligible relationship
9 with influence ($r \approx 0$), while specific tech stack choices (e.g., Java, TypeScript) are strong predictors of
10 popularity. Furthermore, we find that engagement levels do not significantly impact issue resolution time
11 ($p > 0.05$).
12

13 CCS Concepts: • Human-centered computing → Collaborative and social computing; • Software and
14 its engineering → Software maintenance tools.
15

16 Additional Key Words and Phrases: AIDev, software engineering, developer influence, issue resolution, GitHub
17 analysis
18

19 1 Introduction

20 The rapid growth of Artificial Intelligence (AI) development has fostered a massive community
21 of developers on platforms like GitHub. Understanding what drives influence and efficiency in
22 this specific domain is crucial for both individual developers seeking to grow their impact and
23 organizations aiming to optimize their workflows. This project analyzes the AI Dev dataset, focusing
24 on user activity, repository metadata, and issue tracking [?]. The primary objective is to identify
25 the factors that drive developer influence (popularity) and efficiency (issue resolution) within this
26 specific domain.
27

28 2 Research Questions

29 In this study, we investigate the following three research questions (RQs):
30

- 31 • **RQ1: Contribution Volume vs. Influence.** To what extent does repository creation
32 frequency correlate with follower count, and does this vary by programming language?
- 33 • **RQ2: Predictors of Popularity.** Which developer features (account tenure, language,
34 repository count) are the strongest predictors of a user's follower count?
- 35 • **RQ3: Factors Affecting Issue Resolution Time.** Does higher engagement (comment
36 volume) affect the time-to-resolution (TTR) for issues?
37

38 3 Methodology

39 This section details the data wrangling and statistical methodology used to address the RQs.
40

41 3.1 Data Preprocessing and Wrangling

42 We utilized Python (Pandas) to clean and merge the AI Dev dataset. Key preprocessing steps included:
43

- 44 (1) **Complex Merging:** We performed an inner join between the Users and Repositories tables
45 by extracting the username from the repository full_name string (splitting “owner/repo”)
46 to match the login column.
47

48 Authors' Contact Information: Zihuan Lui, University of British Columbia, Kelowna, Canada; Avalvir Kaur Sekhon, University
49 of British Columbia, Kelowna, Canada.

- 50 (2) **Issue-Comment Linkage:** We merged the Issues table with aggregated data from the
 51 Comments table using id and pr_id as foreign keys to calculate engagement metrics.
 52 (3) **Time Standardization:** All timestamp columns (created_at, closed_at) were converted
 53 to UTC. We derived account_age_days and calculated Time-to-Resolution (TTR_hours)
 54 by subtracting issue creation time from close time [?].
 55 (4) **Handling Missing Data:** We imputed missing values for main_language with “Unknown”
 56 and applied One-Hot Encoding to categorical language variables for regression analysis.
- 57

58 **3.2 Analysis Approach**

- 59 • **RQ1:** We aggregated user repository counts and correlated them with follower counts. We
 60 applied Spearman’s rank correlation to handle the non-normal distribution of follower
 61 counts, segmented by the top 5 languages.
 62 • **RQ2:** We constructed a Linear Regression model using features such as account_age,
 63 repo_count, and encoded language. We interpreted coefficients to determine feature im-
 64 portance.
 65 • **RQ3:** We binned issues into “High” and “Low” engagement based on the median comment
 66 count. A Two-sample T-test was used to determine if there is a statistically significant
 67 difference in TTR between these groups.
- 68

69 **4 Results**

70 The following subsections present the statistical results of our analysis.

71

72 **4.1 RQ1: Contribution Volume vs. Influence**

73 Our analysis using Spearman’s rank correlation reveals that the volume of contributions (repository
 74 count) has a negligible relationship with user influence (followers) across all major languages. The
 75 correlation coefficients are consistently near zero:
 76

- 77 • **Python:** 0.0252
 78 • **Java:** 0.0340
 79 • **JavaScript:** -0.0111
 80 • **HTML:** -0.0152
- 81

82 **4.2 RQ2: Predictors of Popularity**

83 The Linear Regression model identified distinct drivers for follower counts:

84

- 85 • **Positive Drivers:** The strongest predictors were specific languages: **Java (+4.27)** and
 86 **TypeScript (+3.21)**. Repository count had a moderate positive impact (+3.72).
 87 • **Negative Drivers:** Users primarily associated with PHP (-13.7) and HTML (-6.3) tended to
 88 have fewer followers.
 89 • **Neutral Factors:** Account tenure (account_age_days) had a coefficient near zero (+0.014).
- 90

91 **4.3 RQ3: Factors Affecting Issue Resolution Time**

92 We analyzed whether higher engagement (median comment count = 5.0) impacted resolution time.
 93 The Two-sample T-test yielded a p-value of ~ 0.30 ($t \approx -1.04$). Since $p > 0.05$, we fail to reject the
 94 null hypothesis; there is no statistically significant difference in resolution time between high and
 95 low engagement issues.

96

99 5 Interpretation of Results

100 5.1 Quantity vs. Quality

101 The findings from RQ1 suggest that developers cannot simply increase their influence by creating a
102 large volume of repositories. The lack of correlation implies that the community values the quality
103 or utility of a project over sheer quantity.
104

105 5.2 Tech Stack Influence

106 RQ2 results indicate that tech stack choice outweighs account longevity. The strong positive
107 coefficients for Java and TypeScript suggest these ecosystems currently offer higher visibility.
108 Conversely, the neutral impact of account age suggests that newer developers can gain influence
109 quickly if they contribute to the right ecosystems.
110

111 5.3 Efficiency Dynamics

112 The RQ3 results suggest that discussion volume does not inherently slow down or speed up the
113 resolution process. Future work will investigate if “Repository Stars” (project popularity) have a
114 more significant effect on TTR than comment engagement.
115

116 6 Project Resources

- GitHub Repository: <https://github.com/RoxyLiu66/data-wrangling-group-8.git>
- Dataset Source: AIDev Dataset

120 7 Remaining Work

121 Milestone 3 will address: Do merge outcomes differ between Agentic and Human PRs after con-
122 trolling for PR size and reviewer activity? Planned analyses include logistic regression, bootstrap
123 confidence intervals, and per-language comparisons.
124

125 8 GenAI Usage Statement

126 We utilized Generative AI (ChatGPT) to assist in debugging Python syntax for the data merging
127 functions and to refine the wording of the methodology section. All code logic and statistical
128 interpretations were verified manually by the team.
129

130 Acknowledgments

131 We would like to thank the course instructors for their guidance on DATA 542.
132

133

134

135

136

137

138

139

140

141

142

143

144

145

146

147