

# 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, contrary to initial assumptions, we find that engagement levels significantly impact  
11 issue resolution time ( $p < 0.001$ ), with engaged issues being resolved faster than those with no activity.  
12

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

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

## 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 3 Methodology

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

### 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 (1) **Complex Merging:** We performed an inner join between the Users and Repositories tables  
44 by extracting the username from the repository full\_name string (splitting "owner/repo")  
45 to match the login column.

46  
47 Authors' Contact Information: Zihuan Lui, University of British Columbia, Kelowna, Canada; Avalvir Kaur Sekhon, University  
48 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. Crucially, we imputed missing comment  
 52 counts with 0 to accurately reflect issues with no engagement.  
 53 (3) **Time Standardization:** All timestamp columns (created\_at, closed\_at) were converted  
 54 to UTC. We calculated Time-to-Resolution (TTR\_hours) by subtracting issue creation time  
 55 from close time, filtering out data errors (negative durations) [? ].  
 56 (4) **Handling Missing Data:** We imputed missing values for main\_language with “Unknown”  
 57 and applied One-Hot Encoding to categorical language variables for regression analysis.

### 59 3.2 Analysis Approach

- 60 • **RQ1:** We aggregated user repository counts and correlated them with follower counts. We  
 61 applied Spearman’s rank correlation to handle the non-normal distribution of follower  
 62 counts, segmented by the top 5 languages.  
 63 • **RQ2:** We constructed a Linear Regression model using features such as account\_age,  
 64 repo\_count, and encoded language. We interpreted coefficients to determine feature im-  
 65 portance.  
 66 • **RQ3:** We categorized issues based on engagement (Has Comments vs. No Comments) using  
 67 the median split (Median = 0). A Two-sample T-test (Welch’s t-test) was used to determine  
 68 if there is a statistically significant difference in TTR between these groups.

## 70 4 Results

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

### 73 4.1 RQ1: Contribution Volume vs. Influence

74 Our analysis using Spearman’s rank correlation reveals that the volume of contributions (repository  
 75 count) has a negligible relationship with user influence (followers) across all major languages. The  
 76 correlation coefficients are consistently near zero:  
 77

- 78 • **Python:** 0.0252
- 79 • **Java:** 0.0340
- 80 • **TypeScript:** 0.0331
- 81 • **JavaScript:** -0.0111
- 82 • **HTML:** -0.0152

### 84 4.2 RQ2: Predictors of Popularity

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

- 86 • **Positive Drivers:** The strongest predictors were specific languages: **Java (+4.27)** and  
**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).

### 92 4.3 RQ3: Factors Affecting Issue Resolution Time

94 We analyzed the impact of engagement on resolution time. The median comment count for the  
 95 dataset was 0. Comparing issues with comments (High Engagement) versus those without (Low  
 96 Engagement), the Two-sample T-test yielded a T-statistic of -18.02 and a p-value of  $3.63 \times 10^{-69}$ .  
 97 Since  $p < 0.05$ , we reject the null hypothesis. The negative T-statistic indicates that issues with

99 engagement (comments) have a significantly **lower** time-to-resolution (Mean  $\approx$  22 hours) compared  
100 to unengaged issues (Mean  $\approx$  2675 hours).

## 101 102 **5 Interpretation of Results**

### 103 104 **5.1 Quantity vs. Quality**

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

### 108 109 **5.2 Tech Stack Influence**

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

### 114 115 **5.3 Efficiency Dynamics**

116 The results from RQ3 provide a compelling insight: engagement accelerates resolution. Issues that  
117 attract community discussion are resolved significantly faster than those that remain silent. This  
118 suggests that “noise” (comments) in the AI community often represents active collaboration or  
119 clarification that aids the maintainer, rather than obstruction or debate that delays the fix.

## 120 **6 Project Resources**

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

## 123 124 **7 Remaining Work**

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

## 128 129 **8 GenAI Usage Statement**

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

## 133 134 **Acknowledgments**

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

136

137

138

139

140

141

142

143

144

145

146

147