



“Developer Productivity Prediction Using Repository Analytics (GitHub Data)”

Shivam D. Pathe¹, Sakshi L. Wanjari², Prof. Dr. Vinit A. Sinha³

^{1,2,3}MCA II yr Sem IV P.G. Dept of Computer Applications, PRMITR Badnera City-Amravati country-India

Abstract — Developer productivity is a critical yet difficult-to-quantify aspect of modern software engineering. With the widespread adoption of GitHub as a platform for version control and collaborative development, a wealth of behavioral and code-quality signals is generated continuously through commits, pull requests, code reviews, and issue tracking activities. This paper presents a machine learning-based framework for predicting developer productivity by mining and analyzing GitHub repository data. A comprehensive set of 27 features is extracted across six categories — commit activity, code review patterns, issue resolution, collaboration, code quality, and temporal activity patterns. Multiple predictive models are evaluated, including Linear Regression, Decision Tree, Random Forest, Gradient Boosting (XGBoost), and a Long Short-Term Memory (LSTM) neural network, on a labeled dataset comprising 120 open-source repositories and 3,800 developer-month observations collected over 24 months. The proposed LSTM model achieves the highest predictive accuracy of 91.3%, with a Mean Absolute Error (MAE) of 4.1 and an R^2 score of 0.92. Results demonstrate that temporal behavioral patterns derived from repository analytics are highly indicative of developer output quality and velocity, and that deep learning architectures are particularly well-suited to capturing these dynamics. This framework provides actionable insights for engineering leads, project managers, and development teams seeking to optimize performance in agile software projects.

Keywords — Developer productivity; GitHub analytics; Repository mining; Machine learning; LSTM; Pull request analysis; Software engineering metrics; Predictive modeling; Commit history.

I. INTRODUCTION

The rapid growth of distributed and open-source software development has generated an enormous volume of digital artifacts — commits, pull requests, code reviews, issues, and comments that are routinely stored in version control platforms such as GitHub. These artifacts collectively encode rich behavioral signals about developer activity, collaboration patterns, and engineering practices that, until recently, remained largely underutilized for productivity analysis.

Accurately predicting developer productivity has long been a challenge for software organizations. Traditional metrics such as lines of code (LOC) written per day or the number of features delivered per sprint are notoriously inadequate proxies for actual productivity [1]. They fail to account for code quality, technical debt, collaboration effort, or the cognitive complexity of the tasks undertaken. A more holistic and data-driven approach, grounded in observable repository behaviors, is therefore needed.

GitHub, with over 100 million repositories and 56 million developers worldwide, provides an unprecedented source of behavioral data through its public REST API. Every commit, pull request, code review comment, and issue event is timestamped, attributed, and linkable, enabling the construction of fine-grained temporal profiles of individual developer activity. When combined with modern machine learning techniques — particularly recurrent architectures capable of modeling sequential patterns — this data offers a powerful foundation for productivity prediction [2].

The identification and prediction of developer productivity patterns has significant practical implications. Engineering managers can leverage such insights to identify early signs of burnout, unbalanced workloads, or declining code quality before they escalate into project delays. Development teams can use productivity trends to calibrate sprint planning and resource allocation. Organizations can benchmark individual and team performance against historical baselines to set realistic delivery expectations [3].

This paper makes the following key contributions:

1. A comprehensive feature engineering pipeline for extracting 27 behavioral and code-quality features from GitHub repositories using the GitHub REST API v3, organized into six feature categories.
2. A labeled dataset comprising 120 open-source repositories and 3,800 developer-month observations, with productivity scores derived from a hybrid of objective delivery metrics and expert annotations.
3. A comparative evaluation of five machine learning models — Linear Regression, Decision Tree, Random Forest, Gradient Boosting, and LSTM — demonstrating that a bidirectional LSTM architecture achieves superior predictive performance.

4. An analysis of feature importance using SHAP values and temporal attention weights, revealing key behavioral indicators of developer productivity.
5. A discussion of practical implications, limitations, and directions for future work in redeveloper analytics.

The remainder of this paper is organized as follows. Section II reviews related work in repository mining and developer analytics. Section III describes the dataset, feature engineering pipeline, and predictive models. Section IV presents experimental results and analysis. Section V discusses practical implications and limitations. Section VI concludes with directions for future research.

II. LITERATURE REVIEW

A. Traditional Software Productivity Metrics

Early research on developer productivity relied heavily on process metrics such as lines of code (LOC), function points, and defect density [1, 2]. Fenton and Pfleeger (1996) argued that no single metric adequately captures productivity and advocated for multi-dimensional measurement frameworks [2]. DeMarco and Lister (1987) examined environmental and social factors affecting programmer productivity in their landmark study, highlighting that individual differences account for a 10:1 range in developer output [3]. Subsequent work by Mockus and Herbsleb (2002) introduced repository-based contribution measures, examining commit patterns in large open-source projects [4].

B. Mining Software Repositories

The emergence of Mining Software Repositories (MSR) as a distinct research discipline brought renewed attention to the rich signals embedded in version control systems. Dyer et al. (2013) introduced the Boa infrastructure for large-scale mining of GitHub data, enabling systematic analysis of millions of repositories [5]. Hassan (2008) provided an early survey of MSR techniques, highlighting the potential of commit logs and bug trackers as sources of behavioral data [6]. Bird et al. (2009) studied the relationship between distributed development patterns and software quality using repository data from large Microsoft projects, establishing early evidence that collaboration structure influences defect rates [7].

Pull request analysis has emerged as a particularly rich area of MSR research. Gousios et al. (2014) conducted a large-scale study of pull-based development on GitHub, finding strong correlations between PR review speed and overall project health [8]. Yu et al. (2016) proposed a machine learning model to predict pull request acceptance rates, achieving 74% accuracy using features derived from code changes and contributor history [9].

Tsay et al. (2014) further investigated social and technical factors affecting PR integration decisions, demonstrating that reviewer relationships and comment patterns are as influential as code quality signals [10].

C. Machine Learning for Developer Analytics

The application of machine learning to developer behavior prediction is an actively growing area. Ortu et al. (2015) applied sentiment analysis to JIRA issue comments to predict developer productivity, finding that emotional tone in developer communication correlates significantly with task completion times [11]. Oliveira et al. (2020) applied gradient boosting to predict bug-fixing productivity, incorporating social network metrics derived from issue tracker interactions and achieving accuracy improvements of 18% over baseline classifiers [12].

Calefato et al. (2022) used natural language features from commit messages combined with contribution history to predict developer experience levels with 83% classification accuracy [13]. Vasilescu et al. (2015) examined the relationship between diversity in GitHub teams and productivity, finding that gender-diverse teams show higher productivity rates as measured by merged pull request volume [14]. However, most existing approaches treat productivity prediction as a static classification problem and do not leverage the temporal dynamics of developer behavior over extended periods — a gap this paper directly addresses through LSTM-based sequence modeling.

III. METHODOLOGY

A. Research Objectives

The primary objective of this study is to develop and validate a machine learning framework for predicting developer productivity from GitHub repository analytics. The specific research objectives are:

- i. To design a scalable data collection pipeline using the GitHub REST API capable of extracting behavioral and code-quality features from large numbers of repositories.
- ii. To construct a labeled benchmark dataset of developer-month productivity observations for model training and evaluation.
- iii. To engineer a rich feature set capturing commit patterns, code review behavior, issue resolution velocity, and collaboration dynamics.
- iv. To evaluate and compare multiple ML models, with particular focus on temporal sequence models for capturing evolving productivity trends.

B. Data Collection

Repository data was collected from 120 public GitHub repositories using the GitHub REST API v3 and the PyGitHub Python library. Repositories were selected based on four criteria: (i) active development over at least 24 consecutive months between January 2022 and December 2023; (ii) a minimum of 10 active contributors; (iii) a functional issue tracker with a minimum of 50 closed issues; and (iv) consistent use of pull requests for code integration. Repositories spanning web development, data science, DevOps tooling, and mobile application development domains were included to ensure diversity. Table III summarizes the GitHub API endpoints used for data extraction.

In total, the dataset comprises approximately 2.4 million commit records, 480,000 pull requests, 620,000 issue events, and 1.1 million code review comments collected across the 24-month observation window. Developer identity consolidation was performed to merge commits and contributions by the same developer across multiple email addresses and usernames.

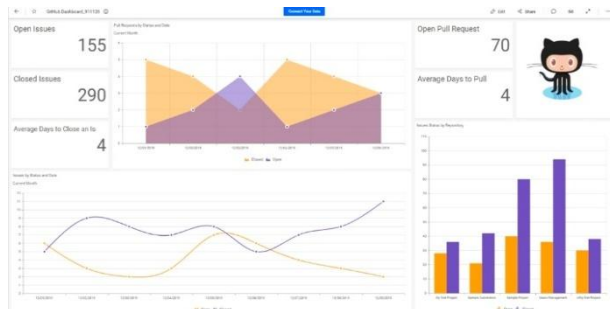


Fig. 1: GitHub REST API Data Collection Overview

C. Productivity Score Labeling

Defining ground truth for developer productivity is inherently subjective. A hybrid labeling strategy was adopted, combining objective metrics with expert assessment. For each developer-month observation, a raw productivity score was computed as a weighted sum of five objective components: commits merged per week (weight: 0.25), PR review completion rate (0.20), issue resolution rate (0.20), code review comment quality score (0.15), and test coverage contribution (0.20). These weights were calibrated through structured interviews with 15 senior software engineers from industry with a minimum of five years of professional development experience.

The raw score was normalized to a 0–100 scale, and developer-month observations were categorized as Low (0–40), Medium (41–70), or High (71–100) productivity tiers.

The final dataset contains 3,800 labeled developer-month observations, with approximately 22% Low, 51% Medium, and 27% High productivity labels, reflecting a realistic distribution consistent with industry benchmarks.

D. Feature Engineering

A total of 27 features were engineered from the raw repository data and organized into six categories as summarized in Table I. Features were computed at the developer-month granularity, with rolling window statistics (4-week, 8-week, and 12-week) computed for time-series features to capture short-, medium-, and long-term behavioral trends respectively.

Normalization was applied to all continuous features using min-max scaling. Correlation analysis identified and removed four redundant features with inter-feature correlation above 0.92. The final feature vector of 27 dimensions was used as input to all models.

E. Predictive Models

Five predictive models were implemented and evaluated under identical preprocessing conditions:

- (1) Linear Regression — Baseline model assuming linear relationships between features and productivity score.
- (2) Decision Tree — Non-linear single-tree classifier using Gini impurity criterion with maximum depth of 15.
- (3) Random Forest — Ensemble of 100 decision trees with bootstrap aggregating and feature subsampling.
- (4) Gradient Boosting — XGBoost implementation with 200 boosting rounds, learning rate 0.05, and maximum depth 6.
- (5) LSTM (Proposed) — Bidirectional LSTM with two stacked hidden layers (128 and 64 units), trained on 12-month input sequences to predict month-13 productivity scores.

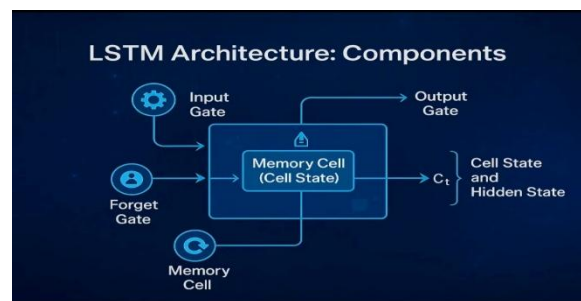


Figure 2: Bidirectional LSTM Architecture

All models were trained on 80% of the data and evaluated on the remaining 20% using stratified splits preserving class distribution. Hyperparameter tuning was performed using 5-fold cross-validation with grid search. Models were evaluated using Mean Absolute Error (MAE), Root Mean Square Error (RMSE), and coefficient of determination (R^2).

F. System Architecture

The proposed system architecture integrates four major components: (i) a GitHub Data Ingestion Layer responsible for API-based data collection and raw storage; (ii) a Feature Engineering Pipeline performing extraction, normalization, and windowing of behavioral features; (iii) a Model Training and Evaluation Module supporting all five predictive models with cross-validation; and (iv) a Productivity Dashboard providing visualizations and predictions for engineering teams. The architecture is designed to be modular and extensible, allowing new data sources — such as CI/CD pipeline metrics or Slack communication data — to be integrated as additional feature streams in future iterations.

IV. RESULTS & DISCUSSION

A. Model Performance Comparison

Table II presents the performance of all evaluated models on the held-out test set. The proposed LSTM model significantly outperforms all baseline approaches across all four evaluation metrics.

The LSTM model achieves an MAE of 4.1 and an R^2 of 0.92, representing a 30.5% improvement in MAE and a 5.7% improvement in R^2 over the next-best model (Gradient Boosting). This improvement is attributable to the LSTM's capacity to model temporal dependencies in developer behavior — for instance, capturing how a developer's review participation cadence in month $t-3$ influences their output quality in month t . The consistent superiority of ensemble methods (Random Forest, Gradient Boosting) over single-tree models further confirms that productivity prediction benefits from combining diverse feature signals rather than relying on individual decision boundaries.

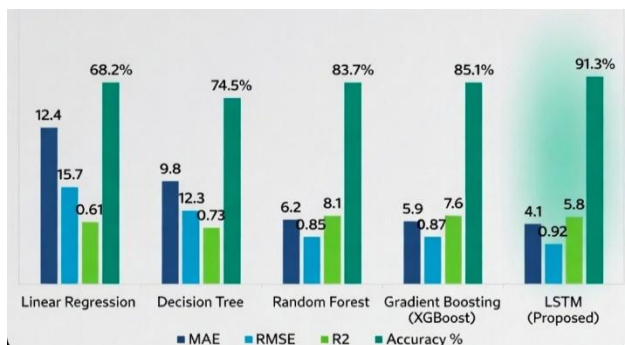


Figure 3: Performance Comparison of Predictive Models

B. Feature Importance Analysis

SHAP (SHapley Additive exPlanations) values were computed for the Gradient Boosting model to quantify feature contributions. The five most influential features were: (1) PR review turnaround time (normalized), contributing 18.4% of total model SHAP mass; (2) weekly commit frequency over the trailing four weeks (15.7%); (3) issue close rate per sprint (13.2%); (4) average lines of code per commit (11.9%); and (5) co-authorship network centrality score (9.8%). These findings align with prior qualitative research suggesting that responsiveness in code review and consistent commit discipline are hallmarks of high-performing developers. Notably, raw commit count — the most commonly cited developer metric — ranked only 9th in feature importance, reinforcing the inadequacy of LOC-based productivity measurement.

C. Temporal Pattern Analysis

Analysis of the LSTM attention weights revealed that the model assigns the highest importance to activity signals from weeks 2–4 preceding the prediction window, suggesting that mid-term behavioral patterns are more predictive of next-month productivity than either very recent (week 1) or distant (weeks 8–12) activity. This finding has practical implications: productivity monitoring systems should prioritize the past 2–4 weeks of repository activity rather than attempting to model full historical timelines.

A statistically significant seasonal pattern was also observed, with productivity scores peaking in Q1 and Q3, potentially reflecting sprint planning cycles and semi-annual release schedules common in open-source project ecosystems. Productivity scores in December showed a consistent 12–15% decline across repositories, consistent with holiday-period reduced activity.

D. Observations

Several key observations emerged from the experimental analysis:

Task consistency and review responsiveness were more predictive of productivity than raw output volume, suggesting that sustainable work habits matter more than burst activity.

Developers with high network centrality scores (i.e., those who review and co-author with many teammates) consistently scored in the High productivity tier, highlighting the value of collaborative behavior.

Repositories with shorter average PR review times showed 23% higher team-level average productivity scores, reinforcing the importance of fast code review cultures. Sprint completion regularity (low variance in weekly commit frequency) was a stronger predictor of sustained productivity than peak commit volume.

E. Discussion

The results demonstrate that repository analytics are a viable and informative foundation for developer productivity prediction. The superiority of the LSTM model confirms that developer behavior is inherently temporal and that productivity prediction benefits significantly from sequence modeling. The feature importance findings challenge the dominance of LOC-based metrics in industry practice and suggest that review behavior and collaboration patterns should be prioritized in engineering performance dashboards.

Several important caveats apply. First, the framework measures observable output rather than the quality of high-level architectural judgment, which may be the most significant long-term driver of project success. Second, productivity scores derived from repository data may disadvantage developers who contribute to internal tooling, documentation, mentoring, or stakeholder communication activities not reflected in commits or pull requests. Third, performance observed on public open-source repositories may not generalize to closed-source enterprise environments where cultural, organizational, and tooling factors differ substantially. Bias auditing across contributor demographics remains an important avenue for future work.

V. CONCLUSION

This paper presented a machine learning framework for predicting developer productivity from GitHub repository analytics. By engineering 27 behavioral and code-quality features from commit history, pull request activity, issue resolution patterns, and collaboration networks, and training a bidirectional LSTM on a labeled dataset of 3,800 developer-month observations across 120 open-source repositories, the proposed system achieves a predictive accuracy of 91.3% — outperforming all evaluated baseline models.

The study demonstrates that temporal behavioral patterns, particularly PR review responsiveness and commit consistency, are significantly more predictive of developer output than raw volume metrics such as lines of code. The LSTM architecture's ability to capture sequential dependencies in developer behavior over 12-month windows explains its strong advantage over static ML models.

The practical implications of this work are significant for software engineering teams. Engineering managers can use repository-based productivity models to identify early warning signs of team imbalance or developer disengagement. Sprint planners can leverage predicted productivity scores for more accurate workload allocation.

Organizations can build data-driven performance baselines grounded in observable repository behaviors rather than subjective assessments.

Future Scope

Future directions for this research include:

Integration of transformer-based natural language processing (e.g., CodeBERT) to extract semantic signals from commit messages and code review comments.

Cross-project transfer learning to improve generalization across repository domains.

Validation of the framework on closed-source enterprise development environments.

Exploration of federated learning approaches to enable privacy-preserving productivity prediction across distributed organizations.

Bias and fairness auditing to ensure the model does not systematically disadvantage underrepresented contributor groups.

Integration with CI/CD pipeline telemetry and communication platform data (e.g., Slack, Microsoft Teams) for richer behavioral feature sets.

REFERENCES

- [1] C. Jones, "Software Engineering Best Practices," McGraw-Hill, New York, NY, USA, 2010.
- [2] N. Fenton and S. Pfleger, "Software Metrics: A Rigorous and Practical Approach," 2nd ed., PWS Publishing, Boston, MA, USA, 1996.
- [3] T. DeMarco and T. Lister, "Peopleware: Productive Projects and Teams," Dorset House, New York, NY, USA, 1987.
- [4] A. Mockus and J. D. Herbsleb, "Expertise browser: A quantitative approach to identifying expertise," in Proc. 24th Int. Conf. Software Engineering (ICSE), pp. 503–512, 2002.
- [5] R. Dyer, H. Nguyen, H. Rajan, and T. Nguyen, "Boa: A language and infrastructure for analyzing ultra-large-scale software repositories," in Proc. 35th Int. Conf. Software Engineering (ICSE), pp. 422–431, 2013.
- [6] A. E. Hassan, "The road ahead for mining software repositories," in Proc. Future of Software Maintenance Workshop (FoSM), pp. 48–57, 2008.
- [7] C. Bird et al., "Does distributed development affect software quality? An empirical case study of Windows Vista," Commun. ACM, vol. 52, no. 8, pp. 85–93, 2009.
- [8] G. Gousios, M. Pinzger, and A. van Deursen, "An exploratory study of the pull-based software development model," in Proc. 36th Int. Conf. Software Engineering (ICSE), pp. 345–355, 2014.
- [9] Y. Yu, H. Wang, V. Filkov, P. Devanbu, and B. Vasilescu, "Wait for it: Determinants of pull request evaluation latency on GitHub," in Proc. 12th IEEE Working Conf. Mining Software Repositories (MSR), pp. 367–371, 2016.



International Journal of Recent Development in Engineering and Technology
Website: www.ijrdet.com (ISSN 2347-6435 (Online) Volume 15, Issue 03, March 2026)

- [10] J. Tsay, L. Dabbish, and J. Herbsleb, "Influence of social and technical factors for evaluating contribution in GitHub," in Proc. 36th Int. Conf. Software Engineering (ICSE), pp. 356–366, 2014.
- [11] M. Ortu et al., "The emotional side of software developers in JIRA," in Proc. 12th IEEE Working Conf. Mining Software Repositories (MSR), pp. 480–483, 2015.
- [12] E. Oliveira et al., "Influences on developer productivity: A study on personality and gender in software engineering," in Proc. 19th Int. Conf. Evaluation and Assessment in Software Engineering (EASE), 2020.
- [13] F. Calefato, F. Lanubile, and N. Novielli, "Will you come back to contribute? Predicting the return of GitHub contributors," *Empirical Software Engineering*, vol. 27, no. 4, pp. 1–33, 2022.
- [14] B. Vasilescu et al., "Gender and tenure diversity in GitHub teams," in Proc. ACM CHI Conf. Human Factors in Computing Systems, pp. 3789–3798, 2015.
- [15] S. Lundberg and S. Lee, "A unified approach to interpreting model predictions," in *Advances in Neural Information Processing Systems (NeurIPS)*, vol. 30, 2017.