

Tapajit Dey

· COMPUTER SCIENCE — EMPIRICAL SOFTWARE ENGINEERING — INNERSOURCE — MINING SOFTWARE REPOSITORIES

☎ (+353) 89 9528287 | @tapajit.dey [at] lero.ie, tdey2 [at] vols.utk.edu | 🌐 tapjey.github.io | 📄 tapjey | 📄 tapajit-dey | 🎓 Google Scholar - Tapajit Dey

📄 Summary

- Postdoctoral Researcher at Lero - the Science Foundation Ireland Research Centre for Software, working on InnerSource software development methodology adoption.
- Experienced in Software Repository Mining, Quantitative Research methods including various Supervised and Unsupervised Machine Learning techniques, Empirical Software Engineering, and Open Source Software.

📁 Work Experience

Lero, the Science Foundation Ireland Research Centre for Software

Limerick, Ireland

POSTDOCTORAL RESEARCHER

September 2020 - Present

- Responsibility - To undertake research on Innersource software development/Corporate Software Engineering Capability for Huawei as part of the TREES Research Centre.

University of Tennessee, Knoxville

Knoxville, TN, USA

GRADUATE RESEARCH AND TEACHING ASSISTANT

August 2015 - August 2020

- TA Courses: Fundamentals of Digital Archeology, Programming Languages and Systems, Operating Systems

IBM India — India Software Lab : SRDC

Bangalore, India

STAFF RESEARCH ENGINEER

June 2012 - June 2015

- **Primary Responsibility** — Key member of the team in charge of Quality Assurance of IBM's 200mm Process Development Kits (PDK), containing device models (primarily MOSFETs) and support components for chip design and layout
- Automated the test procedures using **Python** and created an end-to-end tool for running simulations, parsing the result, and creating a PDF report; that enhanced test coverage by ~**300x**, reduced turn-around-time by **75%**, improved efficiency by **50%**, and reduced customer issues drastically
- Test methodology: "An Integrated Approach to Process Design Kit (PDK) Verification" accepted for **Patent Publication** (Ref. No. IN820131127)

Python Quality Assurance Automation Tool Design

IBM India — India Software Lab : SRDC

Bangalore, India

INTERNSHIP

May 2011 - July 2011

- Designed a Python based GUI for simulating and plotting the results of simulation for SOI (Silicon On Insulator) devices

🎓 Education

University of Tennessee, Knoxville (UTK)

Knoxville, TN, USA

PH.D. IN COMPUTER SCIENCE, WITH DR. AUDRIS MOCKUS

Aug. 2015 - Aug. 2020

- *Thesis Title* — Modeling User-Affected Software Properties for Open Source Software Supply Chains
- *Research Focus* — Software Repository Mining, Data Analytics, Empirical Software Engineering

Indian Institute of Technology (IIT) Kharagpur

Kharagpur, WB, India

DUAL DEGREE (B.TECH.(HONS.) + M.TECH.) IN ELECTRONICS AND ELECTRICAL COMMUNICATION ENGINEERING

July 2007 - May 2012

- *Specialization* — Microelectronics and VLSI Design
- Overall GPA — 8.67/10.00

📄 Skills

PROGRAMMING LANGUAGES (ACCORDING TO LEVEL OF PROFICIENCY)

Python ●●●●● R ●●●●● Shell Script(Bash) ●●●○○○
MySQL ●●●○○ Perl ●●○○○ JavaScript/HTML/C/C++ ●○○○○

SOFTWARE TOOLS AND MACHINE LEARNING TECHNIQUES

Development Frameworks	R-Studio, Jupyter (IPython) Notebook, Vim, Sublime
Other Tools	MongoDB, Git, Docker, LaTeX
Supervised Learning	Linear Regression (LM, GLM), Random Forest, Neural Network (Deep Learning), SVM
Unsupervised Learning	Clustering - K-means, Hierarchical etc.; Dimension Reduction - PCA, SVD; Association Rule
Bayesian Networks (BN)	BN structure search, BN inference (Causal and Non-Causal)
Research Specific Tools	World of Code tool, GHTorrent, Google Big Query
Languages Spoken	English, Bengali, Hindi

≡ Relevant Research Projects

Effects of Software Usage on Perceived Quality of a Software

- Proposed an usage-independent *actionable* and *easy-to-use measure* of software quality
- Used **Bayesian Networks** and **Random Forest** to model the interrelationship between software faults, complexity, and usage
- Developed the model by studying a proprietary mobile software and verified the external validity of the concept for popular NPM packages
- **Involvement:** Methodology development and testing, Authored 1 Conference and 1 Journal paper.

Software Usage Software Quality Bayesian Network Random Forest R

Study of the NPM Ecosystem: User Participation and Package Dependencies (Ongoing)

- Created Software Supply Chain for NPM Ecosystem by analyzing dependencies of packages
- Explored User participation to NPM packages inside and outside of the users' individual supply chains, i.e. packages the user's projects depend on directly or indirectly (dataset consisted of **1,376,946** issues and pull-requests(PR) created for **4433** NPM packages and **272,142** issue creators)
- Classified the Users based on their participation patterns using C-means clustering algorithm
- Investigated the factors affecting Pull Request Quality (chance of the Pull Request being accepted within a month) using Random Forest models
- Replicated a study that explored a similar question by investigating a larger dataset (**483,988** Pull Requests from **4218** popular NPM packages) and improved the performance by adding additional factors, achieving an *AUC-ROC value of 0.95*, **6%** higher than the previous study
- **Involvement:** Methodology Development & Testing, Authored 2 conference papers and 1 book chapter.

Software Supply Chain Software Dependency NPM Ecosystem Issues Pull Requests Clustering R Python

Study of Bots that Commit Code (Ongoing)

- Proposed a systematic approach to detect bots that commit code in Git repositories using committer name, commit messages, files modified by the commit, and projects associated with the commits.
- Characterized bots based on the commits, identified frequently modified file types.
- Collaborated with other Graduate and Undergraduate Students, an industry professional from GitHub, and a professor from Carnegie Mellon University.
- **Involvement:** Methodology Development & Testing, Leading Collaboration for Project, Supervising Other Graduate and undergraduate students during project, Authored 1 conference and 1 workshop paper (invited).

Bot Detection Code Commits Text Analysis Random Forest Predictive Modeling Python R

Representing Software Developers' and Projects' specific expertise using Vector Embedding, Project and Developer aliasing (Ongoing)

- Used Libraries imported in code files edited/authored by developers/projects to measure their expertise.
- Used LSI and Doc2Vec for representing their expertise for over a skill space.
- Developed novel methods for linking developer IDs using their name, email, and GitHub IDs at scale (40M developers)
- Used community detection algorithm for detecting project forks.
- **Involvement:** Development of methodology, Supervision of undergraduate and graduate students for testing the approaches, Authored 1 conference paper.

LSI Doc2Vec Expertise Python Open Source Software

Honors & Awards

- 2019 **Best Paper Award**, Proceedings of the 15th International Conference on Predictive Models and Data Analytics in Software Engineering (PROMISE - 2019) *Recife, Brazil*
- 2017 **Best Graduate Research Assistant**, University of Tennessee, Knoxville — College of Engineering (Award Amount: \$ 2500) *Knoxville, TN, USA*
- 2011 **Finalist**, Best B.Tech Project, for designing an ultra low power watchdog timer — Department of Electronics and Electrical Communication Engineering, Indian Institute of Technology, Kharagpur *Kharagpur, WB, India*
- 2005 **Ranked 4th**, Regional Mathematical Olympiad, (Selection test for International Mathematical Olympiad) *India*

Publications: Tapajit Dey

BOOK CHAPTERS

- [1] Sadika Amreen, Bogdan Bichescu, Randy Bradley, **Tapajit Dey**, Yuxing Ma, Audris Mockus, Sara Mousavi, and Russell Zaretski. “A Methodology for Measuring FLOSS Ecosystems.” In “Towards Engineering Free/Libre Open Source Software (FLOSS) Ecosystems for Impact and Sustainability”, pp. 1-29. Springer, Singapore, 2019. [LINK](#)

JOURNAL ARTICLES

- [1] Yuxing Ma, **Tapajit Dey**, Chris Bogart, Sadika Amreen, Marat Valiev, Adam Tutko, David Kennard, Russell Zaretski, Audris Mockus “World of Code: Enabling a Research Workflow for Mining and Analyzing the Universe of Open Source VCS data”, Accepted in the Empirical Software Engineering journal.
- [2] Andrey Krutauz, **Tapajit Dey**, Peter Rigby, and Audris Mockus. “Do Code Review Measures Explain the Incidence of Post-Release Defects?” Empirical Software Engineering 25, no.5 (2020): 3323–3356. [LINK](#). Impact Factor: 4.457 (2018)
- [3] **Tapajit Dey** and Audris Mockus. “Deriving a Usage-Independent Software Quality Metric.” Empirical Software Engineering (2020): 1-46. [LINK](#). Impact Factor: 4.457 (2018)

CONFERENCE ARTICLES

- [1] **Tapajit Dey**, Andrey Karnauch, and Audris Mockus “Skill Spaces: Representation of Expertise in Open Source Software.” ACCEPTED in ICSE-2021. [LINK](#)
- [2] **Tapajit Dey** and Audris Mockus “Effect of Technical and Social Factors on Pull Request Quality for the NPM Ecosystem.” In Proceedings of the 14th ACM/IEEE International Symposium on Empirical Software Engineering and Measurement (ESEM), pp. 1-11. 2020. [LINK](#)
- [3] **Tapajit Dey**, Sara Mousavi, Eduardo Ponce, Tanner Fry, Bogdan Vasilescu, Anna Filippova, and Audris Mockus “Detecting and Characterizing Bots that Commit Code.” In Proceedings of the 17th International Conference on Mining Software Repositories (pp. 209–219). Association for Computing Machinery, 2020. [LINK](#)
- [4] Tanner Fry, **Tapajit Dey**, Andrey Karnauch, and Audris Mockus “A Dataset and an Approach for Identity Resolution of 38 Million Author IDs extracted from 2B Git Commits.” In Proceedings of the 17th International Conference on Mining Software Repositories (pp. 518–522). Association for Computing Machinery, 2020. [LINK](#)
- [5] **Tapajit Dey**, Yuxing Ma, and Audris Mockus. “Patterns of effort contribution and demand and user classification based on participation patterns in npm ecosystem.” In Proceedings of the Fifteenth International Conference on Predictive Models and Data Analytics in Software Engineering, pp. 36-45. ACM, 2019. [LINK](#)

- [6] **Tapajit Dey and Audris Mockus.** “Modeling Relationship between Post-Release Faults and Usage in Mobile Software.” In Proceedings of the 14th International Conference on Predictive Models and Data Analytics in Software Engineering, pp. 56-65. ACM, 2018. [LINK](#)
- [7] **Tapajit Dey and Audris Mockus.** “Are software dependency supply chain metrics useful in predicting change of popularity of npm packages?.” In Proceedings of the 14th International Conference on Predictive Models and Data Analytics in Software Engineering, pp. 66-69. ACM, 2018. [LINK](#)
- [8] **Tapajit Dey, Jacob Logan Massengill, and Audris Mockus.** “Analysis of popularity of game mods: A case study.” In Proceedings of the 2016 Annual Symposium on Computer-Human Interaction in Play Companion Extended Abstracts, pp. 133-139. ACM, 2016. [LINK](#)

WORKSHOP ARTICLES

- [1] **Tapajit Dey, Bogdan Vasilescu, and Audris Mockus** “An Exploratory Study of Bot Commits.” In Proceedings of the IEEE/ACM 42nd International Conference on Software Engineering Workshops (pp. 61–65). Association for Computing Machinery, 2020. [LINK](#)
- [2] **Yuxing Ma, Tapajit Dey, and Audris Mockus.** “Modularizing global variable in climate simulation software: position paper.” In Proceedings of the International Workshop on Software Engineering for Science, pp. 8-11. ACM, 2016. [LINK](#)

PRE-PRINT ARTICLES (NOT PEER-REVIEWED)

- [1] **Tapajit Dey and Audris Mockus.** “A Matching Based Theoretical Framework for Estimating Probability of Causation.” arXiv preprint arXiv:1808.04139 (2018). — Proposing a novel way to estimate probability of causation using statistical matching for experimental and observational scenarios. [LINK](#)
- [2] **Yuxing Ma, Tapajit Dey, Jarred M. Smith, Nathan Wilder, and Audris Mockus.** “Crowdsourcing the discovery of software repositories in an educational environment.” PeerJ Preprints 4 (2016): e2551v1. [LINK](#)

WORKS UNDER REVIEW/IN PROGRESS

- [1]

Citations: — See List of Publications from the [Google Scholar](#) Page for citation information

Other Activity

- Participated in multiple Hackathons and Webinars about the World of Code tool.
- Member of ACM.
- Shared Multiple Datasets related to published work for ease of community access and enabling replication:
 - <https://zenodo.org/record/3701819>
 - <https://zenodo.org/record/3648702>
 - <https://zenodo.org/record/3653283>
 - <https://zenodo.org/record/3700859>
 - <https://zenodo.org/record/3699665>
 - <https://zenodo.org/record/3694401>