

# FlightRecorder: Tracing Fine-Grained Development with AI Assistants

Tobias Dick  
Saarland University  
Saarbrücken, Germany  
tdick@cs.uni-saarland.de

Benedict Bliem  
Saarland University  
Saarbrücken, Germany  
bebl00001@stud.uni-saarland.de

Kallistos Weis  
Saarland University  
Saarbrücken, Germany  
kallistos@cs.uni-saarland.de

Alisa Welter  
Saarland University  
Saarbrücken, Germany  
welter@cs.uni-saarland.de

Sven Apel  
Saarland University  
Saarbrücken, Germany  
apel@cs.uni-saarland.de

## Abstract

Software development is rapidly transforming with the widespread adoption of AI coding assistants such as GITHUB COPILLOT. While research has begun to investigate how these tools change the development process, fundamental questions about how code evolves as developers collaborate with AI agents, particularly at fine-grained levels of interaction, remain open. A key limitation of existing research is reliance on telemetry data that cannot distinguish agent-generated changes from inline completions and human edits, nor track when and where modifications occur with sufficient precision. We present FLIGHTRECORDER, a lightweight VSCode extension that collects fine-grained code evolution data without disrupting the workflows of developers. Our extension integrates three data sources: snapshots of agent-mode edits, inline completions, and chat history, all timestamped for precise attribution. Researchers can export these data in structured formats alongside anonymization and basic visualization functionality. By making such fine-grained evolution data accessible, FLIGHTRECORDER enables researchers to investigate fundamental questions about AI-assisted development. We provide a replication package, including the code for FLIGHTRECORDER and the data used throughout our paper.<sup>1</sup>

### ACM Reference Format:

Tobias Dick, Benedict Bliem, Kallistos Weis, Alisa Welter, and Sven Apel. 2026. FlightRecorder: Tracing Fine-Grained Development with AI Assistants. In . ACM, New York, NY, USA, 5 pages. <https://doi.org/10.1145/nnnnnnn.nnnnnnn>

## 1 Introduction

Software development is undergoing a significant shift with the widespread adoption of AI coding assistants such as GITHUB COPILLOT<sup>2</sup>. These tools augment developers through multiple interaction

<sup>1</sup>We further provide a screencast demo of our tool at <https://youtu.be/bZR6aPROVMU>.

<sup>2</sup><https://github.com/features/copilot>

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from [permissions@acm.org](mailto:permissions@acm.org).  
*Conference'17, Washington, DC, USA*

© 2026 Copyright held by the owner/author(s). Publication rights licensed to ACM.  
ACM ISBN 978-x-xxxx-xxxx-x/YYYY/MM  
<https://doi.org/10.1145/nnnnnnn.nnnnnnn>

```
207 def handle_create(self, name: str, dish_type: Optional[str] = None, 73
208 dietary_restriction: Optional[Any] = None) -> Dict[str, Any]: 74
209     """Create a new empty meal.""" 75
210     if name in self.meals: 76
211         return {"status": "error", "message": "Meal already exists."} 77
212     dietary_restrictions = self.normalize_dietary_restrictions 78
213     (dietary_restriction) 79
214     self.meals[name] = { 80
215         "ingredients": [], 81
216         "dishType": dish_type or self.config.get("dishType"), 82
217         "dietaryRestriction": dietary_restriction 83
218     } 84
219     return {"status": "ok"} 85
```

Figure 1: Edit history visualization with code colored by the edit that introduced it (oldest → newest).

modes, ranging from inline code completions within the editor to agent-driven development, where autonomous agents handle iterative cycles of implementation, testing, and refinement. This range of capability has made AI-assisted development prevalent in industry and has motivated considerable research on LLMs, agents, and coding assistants [3, 5, 11]. Although prior work has examined how AI assistance shapes development processes, human-agent collaboration, code evolution, software quality, and long-term engineering practice remain poorly understood [7].

Most existing studies on how developers use AI coding assistants rely on behavioral telemetry—tracking IDE focus, clipboard activity, and pasting patterns—rather than analyzing code evolution directly [2, 4, 8]. A key obstacle is the absence of fine-grained observational data. The crux is that existing tools cannot reliably distinguish agent-generated edits from inline completions and human edits, nor can they track when and where these occur without disrupting development workflows.

To address this gap, we present FLIGHTRECORDER, a VSCode<sup>3</sup> extension that collects fine-grained code attribution data when developers use GITHUB COPILLOT, without disrupting their workflows. Beyond data collection, FLIGHTRECORDER includes anonymization and visualization functionality (see Figure 2) to support ethical research practices and enable initial exploration of the collected data. Our extension captures three complementary data sources: (1) snapshots of code edits made by GITHUB COPILLOT’s agent mode, (2) accepted inline code completions, and (3) the agent’s chat history.

Each data source is timestamped allowing researchers to correlate modifications across all three sources and reconstruct the

<sup>3</sup><https://code.visualstudio.com>

complete picture of code evolution. Since agents can perform many code edits in rapid succession, FLIGHTRECORDER does not capture every individual edit. Instead, it accumulates edits until the agent has not made code edits within a user-configurable debounce window. At these pauses, FLIGHTRECORDER creates snapshots of the accumulated agent edits. These snapshots, combined with inline completion logs and chat history, are then exported in a structured format that researchers can analyze.

By collecting these data with precise timestamps, researchers can reconstruct and attribute code evolution, identifying which lines were added or modified, when they were edited, and which tool or interaction triggered each edit, and thus can finally investigate the concrete questions that matter to understand AI-assisted development. How do agents and humans work together to structure source code? Do agents make edits linearly, moving sequentially through files, or do they follow conceptual logic and control flow? How much effort do developers invest in refactoring or building upon agent suggestions? Examining these patterns reveals how code evolves when humans and agents work together.

We evaluate FLIGHTRECORDER in a controlled study with 18 participants using GITHUB COPILOT working on an open-ended development task either alone (human-AI) or in pairs (human-human-AI). Our study demonstrates that FLIGHTRECORDER reliably captures meaningful data about code evolution and developer-agent interaction. We further show how to analyze the collected data through FLIGHTRECORDER's built-in edit history visualization and discuss potential applications for research and practice.

## 2 Background and Related Work

AI coding assistants, such as GITHUB COPILOT, have become widely adopted in software engineering practice. Initially, these tools aimed to increase developer productivity by reducing manual coding effort and accelerating routine tasks, for example through real-time code suggestions within an IDE. Recent advances have enabled AI coding assistants to function as autonomous agents capable of executing multi-step software engineering tasks with little to no human intervention. Consequently, a growing body of work has begun to investigate how developers interact with these agents, as well as how such systems influence software development practices and developer behavior [6, 10, 11].

To track how developers interact with AI coding assistants, prior work often relies on screen recordings and manual annotation of coding sessions [1, 10]. Going further, Guglielmi et al. [4] track developer behavior using an IDE plugin, called TAKO, screen recordings, and survey data to compare tasks performed with and without GITHUB COPILOT, investigating how code suggestions affect developer productivity. TAKO captures inline code completions—the primary interaction mode with GITHUB COPILOT at the time of the study—and provides a dedicated analytics view summarizing usage statistics. However, as GITHUB COPILOT has evolved, its inline completion tracking appears to no longer function, and edits made by the agent are not supported. Tang et al. [8] present CODEGRITS, a JetBrains plugin that combines IDE interaction tracking with eye gaze data to study developers' cognitive processes, and apply it to examine how developers validate and repair LLM-generated code [9], but do not track AI-generated edits directly. Closer to

our work, Basha et al. present CODEWATCHER [2], a client-server system with a VS CODE plugin that captures fine-grained developer interaction events such as insertions, deletions, copy-paste actions and window focus. Javahar et al. [6] use CODEWATCHER to study developers using an AI coding assistant, specifically Amazon's CODEWHISPERER, to uncover behavioral patterns during development. Crucially, CODEWATCHER only has access to IDE telemetry and not to the coding assistant itself, and must therefore infer AI-generated edits through heuristics such as the size of an insertion.

## 3 FlightRecorder

FLIGHTRECORDER is implemented as a lightweight VS CODE extension available on the VS CODE Marketplace<sup>4</sup> and is open source<sup>5</sup>. It runs locally on the developer's machine, integrating seamlessly into their existing workflow with GIT as its only dependency.

**Recording AI Interactions.** FLIGHTRECORDER records AI-related events within explicitly delimited recording sessions that the developer can start and stop on demand. It works by continuously monitoring GITHUB COPILOT's log output for a selected workspace, detecting events that correspond to either accepted *inline completions* or applied *agent edits*. For each detected event, FLIGHTRECORDER records a timestamp and the set of affected files, then immediately captures a snapshot of the current file state. Additional event-specific metadata, including the event type, timestamp, and a list of affected files are stored alongside each snapshot in a structured JSON format. Snapshots are persisted as commits in a GIT repository that FLIGHTRECORDER maintains automatically in the background, storing the recorded history in a well-established, portable format and allowing later analyses to build on existing GIT tooling. To avoid inflating the snapshot history during rapid successive edits, events that fall within a configurable debounce window are collapsed into a single commit, with all individual event-metadata kept intact. Beyond event-based snapshots, FLIGHTRECORDER exports two complementary artifacts either on request or when a recording session ends: (1) a *chat session export*, which captures the conversation files that VS CODE stores locally for the chat interface, and (2) the raw *log file* that GITHUB COPILOT writes to disk during the session. Both artifacts are stored verbatim in the repository, with the only transformations applied being the path and identity anonymizations described below.

**Anonymization.** For FLIGHTRECORDER to be usable in research settings, the collected repositories must not leak personal information such as the developer's identity or local file system structure. To this end, FLIGHTRECORDER provides a command that produces an anonymized clone of the collected repository by rewriting its full GIT history using `git-filter-repo`<sup>6</sup>. Author and committer identities are replaced with stable, deterministic placeholders, ensuring that co-author relationships are preserved across the full history. Absolute paths in file contents and commit messages are replaced as well. Users can configure whether to replace only paths that identify the developer's machine or all absolute paths.

<sup>4</sup><https://marketplace.visualstudio.com/items?itemName=SE-SIC.flight-recorder>

<sup>5</sup><https://github.com/se-sic/flight-recorder>

<sup>6</sup><https://github.com/newren/git-filter-repo>

**Edit History Visualization.** To support qualitative analysis of AI-assisted development data, FLIGHTRECORDER provides an interactive edit history visualization that attributes each character of code to the commit that introduced it. Characters are color-coded by their originating commit, and hovering over any region surfaces the contextual metadata recorded at that point, including the event type, timestamp, and affected files. The color scheme can be configured in three modes: two timeline views that position commits relative to either the full repository history or the current file, and an event-type view that distinguishes between inline completions and agent edits. Together, these views allow researchers to trace how AI-generated code evolves across a session, identifying which regions survive intact, which are iteratively revised, and which are discarded shortly after being introduced.

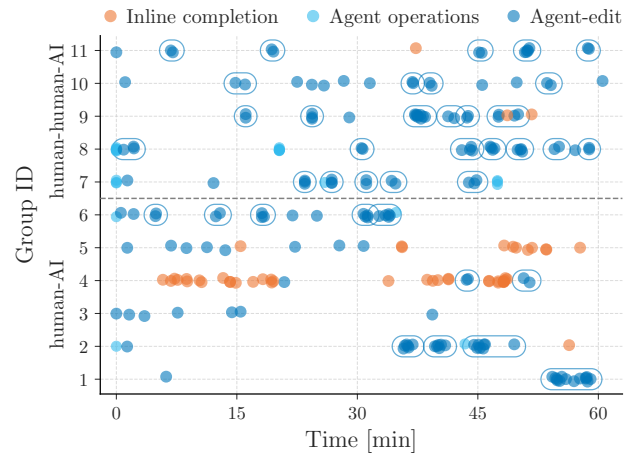
**Scope & Limitations.** FLIGHTRECORDER is designed to capture GITHUB COPILOT activity during active development sessions within a single VSCODE workspace, encompassing both GITHUB COPILOT-induced code changes and chat sessions. All recording is strictly scoped to that workspace; no data from other projects or workspaces is collected, and nothing is written outside the local repository. Since changes are recorded via GIT commits, only accepted inline completions are captured cleanly, as rejected suggestions leave no trace. Agent edits present a subtler challenge: because acceptance is a VSCODE feature rather than a GITHUB COPILOT one, FLIGHTRECORDER cannot intercept it, and instead captures the full edit as soon as it is applied, before the developer has had a chance to review it. A subsequent rejection can therefore only be inferred from later snapshots rather than being attributed to GITHUB COPILOT directly. While FLIGHTRECORDER does not establish a direct link between agent edits and user prompts, each agent-edit snapshot includes the corresponding request identifier, enabling downstream analyses to correlate it with the prompt that triggered the edit.

## 4 Showcase

To evaluate FLIGHTRECORDER, we conduct a controlled empirical study to demonstrate the data it captures and analyses it enables. We report the collected data and key insights enabled by our tool.

**Study Setup.** We recruited 18 computer science students (undergraduate and graduate) from a seminar to participate in our study. Participants were given a specification for a small coding project with two parts: a guided first part with detailed specifications, and an open-ended second part offering extension ideas. Participants either worked individually with GITHUB COPILOT (human-AI,  $n = 8$ ) or in pairs (human-human-AI,  $n = 10$ ). All participants implemented the project in Python over a 60-minute session, using FLIGHTRECORDER to record their activity on their own machines. Prior to the study, all participants provided informed consent.

Refining our study protocol, we conducted a pilot study with several students, which allowed us to improve the project specification and address robustness issues in FLIGHTRECORDER, particularly its handling of GIT misconfigurations and compatibility across log formats from different GITHUB COPILOT versions. Following these refinements, we collected valid data from 16 of 18 participants; one used the command line GITHUB COPILOT rather than VSCODE's integrated version, and another encountered GIT configuration issues mid-session; both were excluded from our analysis.



**Figure 2: Edit timeline across 6 human-AI and 5 human-human-AI groups. Circled clusters denote multiple agent-edits within the same chat response.**

**Workflow Insights.** FLIGHTRECORDER enables us to reconstruct a detailed timeline of AI-related activity for each group, as displayed in Figure 2, with the 60-minute observation window beginning at their first request to GITHUB COPILOT. Recorded events include *inline completions* and *agent operations*, where agent operations encompass file creation, string replacement, and the primary interaction mode *agent edits*. The timeline reveals distinct usage patterns: while agent-edits predominate across most sessions, groups differ notably in their reliance on inline completions and the duration of idle periods between agent interactions.

Interestingly, in the human-human-AI setting, manual editing (i.e., direct code modification or inline completions without invoking the agent) is nearly absent. A manual analysis confirms that human-human-AI pairs worked almost exclusively through GITHUB COPILOT's agent, directing it and interpreting its output rather than writing code themselves. The human-AI condition exhibits greater heterogeneity. Group 6 mirrors the human-human-AI pattern, relying almost exclusively on the agent, while the remaining groups fall into two patterns: groups 4 and 5 alternated between agent use and inline completions, whereas groups 1 and 3 manually edited code without relying on completions. Group 2 represents a technical outlier: although they used the agent between minutes 5 and 35, it operated via the command-line tool rather than VSCODE's edit-tool, creating a gap in recorded data and preventing the capture of individual change acceptances or rejections<sup>7</sup>.

**File-level Inspection.** FLIGHTRECORDER's edit history visualization allows us to inspect how code evolves at the file level for individual participants. Figure 1 shows a function written by group 8, one of the human-human-AI pairs that worked exclusively through the agent. In this function, an initial commit (blue) establishes the core implementation, before later edits (yellow, red) extend the function with additional optional parameters. This progression suggests that the group first created this function when directing the agent to

<sup>7</sup>Group 2 used the weakest available LLM at the time, which likely explains this atypical behavior, especially the failure to follow GITHUB COPILOT's tool-use instructions.

satisfy the guided specification. They later returned to extend it with two separate features from the open-ended part of the project. While this interpretation requires manual inspection to confirm, it illustrates how the edit history visualization can surface structure that is not visible in the session-level timeline: the correspondence between edit events and distinct phases of the task.

## 5 Discussion

The data collected by FLIGHTRECORDER open several avenues for both research and practice. In what follows, we discuss how the collected snapshots can be leveraged to understand agent-assisted development workflows, support developers during interactive coding sessions, and improve code review practices when working with agents. We conclude by outlining promising directions for future extensions of FLIGHTRECORDER.

**Understanding AI-Assisted Development Workflows.** Code evolution is typically tracked through version control systems or direct observation of developers. AI-assisted software development introduces a qualitatively different challenge: agents can generate substantial amounts of code and perform numerous modifications in a short span of time, making it difficult to understand the evolution of the codebase through conventional means. FLIGHTRECORDER addresses this by capturing atomic snapshots of agent edits during a development session, providing fine-grained visibility into how code evolves when developers interact with AI agents. This enables researchers to study AI-assisted development practices in controlled settings and identify recurring interaction patterns, such as which code regions undergo repeated modifications versus which are created once and left intact. Such insights may help understanding how developers leverage agents effectively and which interaction strategies prove most productive.

**Supporting Interactive Development with Agents.** Beyond research applications, FLIGHTRECORDER offers potential to enhance the developer experience during agent-assisted coding. When developers request substantial changes—such as adding features or fixing bugs—agents often generate large changesets that may necessitate iterative refinements or backtracking if the changes prove unsuitable. The snapshot mechanism that FLIGHTRECORDER provides enables developers to easily revert to previous states within a session, offering a form of fine-grained undo that accounts for the scale of agent-generated changes. However, translating this capability to production environments introduces practical challenges. The atomic-level granularity of FLIGHTRECORDER's snapshots, while valuable for detailed analysis, becomes unwieldy in large collaborative codebases, where hundreds of agent edits would generate correspondingly many commits and clutter the project's version control history. To make this approach viable, developers would need to selectively collapse related changes into coherent commits before merging, or work within isolated branches during AI-assisted sessions. Nevertheless, the underlying idea of leveraging fine-grained interaction history to support iterative development merits further exploration for real-world deployment.

**Improving Code Review for Large Changesets.** A critical challenge in AI-assisted development is conducting effective code review when agents introduce hundreds or thousands of lines in a

single session. The edit history captured by FLIGHTRECORDER could inform code review by surfacing heavily-modified code regions as high-priority for more careful review. Distinguishing between code that evolved through iteration (indicating developer uncertainty or agent overcorrection) and code that remained stable from creation gives reviewers additional information about which sections may need increased attention. Such contextual information, derived from the interaction history rather than static analysis, could complement existing review practices.

## 6 Conclusion

We introduce FLIGHTRECORDER, a lightweight VSCODE extension that captures fine-grained code evolution data during AI-assisted development. We show that our tool reliably surfaces behavioral patterns that telemetry-based approaches cannot distinguish, validated through a controlled study with 18 participants.

FLIGHTRECORDER is open source, and we plan to extend it in several directions. Continuously staging human-made edits would cleanly separate human and AI contributions in the repository history. Capturing a wider range of agent tool uses—such as terminal commands or Web searches—would yield more complete records of how agents navigate complex development tasks. We also plan to extend support to additional AI coding assistants such as CLAUDE CODE. We welcome community contributions.

## 7 Data Availability

For reproducibility, we provide a replication package, including the code for FLIGHTRECORDER and the data used throughout our paper at <https://doi.org/10.5281/zenodo.20134366>.

## Acknowledgements

This work has been supported by the European Union as part of the ERC Advanced Grant Brains On Code (101052182) and by the DFG as part of the TRR 248 CPEC (389792660).

## References

- [1] Shradha Barke, Michael B James, and Nadia Polikarpova. 2023. Grounded Copilot: How Programmers Interact with Code-Generating Models. *Proceedings of the ACM on Programming Languages* 7, OOPSLA1 (2023), 85–111. doi:10.1145/3586030
- [2] Manaal Basha, Aimeê M Ribeiro, Jeena Javahar, Cleidson RB De Souza, and Gema Rodríguez-Pérez. 2025. Codewatcher: IDE Telemetry Data Extraction Tool for Understanding Coding Interactions with LLMs. In *2025 IEEE International Conference on Software Maintenance and Evolution*. 935–939. doi:10.1109/icsme64153.2025.00110
- [3] Yihong Dong, Xue Jiang, Jiaru Qian, Tian Wang, Kechi Zhang, Zhi Jin, and Ge Li. 2025. A Survey on Code Generation with LLM-based Agents. *arXiv preprint arXiv:2508.00083* (2025). doi:10.48550/arXiv.2508.00083
- [4] Emanuela Guglielmi, Venera Arnaudova, Gabriele Bavota, Rocco Oliveto, and Simone Scalabrino. 2025. How do Copilot Suggestions Impact Developers' Frustration and Productivity? *arXiv preprint arXiv:2504.06808* (2025). doi:10.48550/arXiv.2504.06808
- [5] Burak Gülmez. 2026. Code Generation with Large Language Models: A Survey from Neural Program Synthesis to Autonomous Software Development. *Applied Intelligence* 56, 6 (2026), 200. doi:10.1007/s10489-026-07230-0
- [6] Jeena Javahar, Tanya Budhrani, Manaal Basha, RB Cleidson de Souza, Ivan Beschastnikh, and Gema Rodríguez-Pérez. 2025. Cracking CodeWhisperer: Analyzing Developers' Interactions and Patterns During Programming Tasks. In *2025 IEEE Symposium on Visual Languages and Human-Centric Computing*. 391–399. doi:10.1109/vl-hcc65237.2025.00051
- [7] Agnia Sergeev, Ilya Zakharov, Ekaterina Koshchenko, and Maliheh Izadi. 2026. Human-AI Experience in Integrated Development Environments: A Systematic

465	Literature Review. <i>Empirical Software Engineering</i> 31, 3 (2026), 55. doi:10.1007/s10664-025-10793-0	<i>Languages and Human-Centric Computing</i> . 40–46. doi:10.1109/vl/hcc60511.2024.00015	523
466			
467	[8] Ningzhi Tang, Junwen An, Meng Chen, Aakash Bansal, Yu Huang, Collin McMillan, and Toby Jia-Jun Li. 2024. CodeGRITS: A Research Toolkit for Developer Behavior and Eye Tracking in IDE. In <i>Proceedings of the 2024 IEEE/ACM 46th international conference on software engineering: Companion proceedings</i> . 119–123. doi:10.1145/3639478.3640037	[10] Alisa Welter, Niklas Schneider, Tobias Dick, Kallistos Weis, Christof Tinnes, Marvin Wyrich, and Sven Apel. 2025. An Empirical Study of Knowledge Transfer in AI Pair Programming. In <i>2025 40th IEEE/ACM International Conference on Automated Software Engineering</i> . 166–177. doi:10.1109/ase63991.2025.00022	524
468			525
469			526
470			527
471	[9] Ningzhi Tang, Meng Chen, Zheng Ning, Aakash Bansal, Yu Huang, Collin McMillan, and Toby Jia-Jun Li. 2024. Developer behaviors in validating and repairing llm-generated code using ide and eye tracking. In <i>2024 IEEE Symposium on Visual</i>	[11] Xin Zhang, Peng Zhang, Yuan Shen, Min Liu, Qiong Wang, Dragan Gašević, and Yizhou Fan. 2024. A Systematic Literature Review of Empirical Research on Applying Generative Artificial Intelligence in Education. <i>Frontiers of Digital Education</i> 1, 3 (2024), 223–245. doi:10.1007/s44366-024-0028-5	528
472			529
473			530
474			531
475			532
476			533
477			534
478			535
479			536
480			537
481			538
482			539
483			540
484			541
485			542
486			543
487			544
488			545
489			546
490			547
491			548
492			549
493			550
494			551
495			552
496			553
497			554
498			555
499			556
500			557
501			558
502			559
503			560
504			561
505			562
506			563
507			564
508			565
509			566
510			567
511			568
512			569
513			570
514			571
515			572
516			573
517			574
518			575
519			576
520			577
521			578
522			579
			580