

Online Appendix for AgenticFlict: A Large-Scale Dataset of Merge Conflicts in AI Coding Agent Pull Requests on GitHub

DANIEL OGENRWOT, University of Nevada Las Vegas, USA

JOHN BUSINGE, University of Nevada Las Vegas, USA

A Online Appendix: Dataset and Artifact Description

A.1 Overview

Software Engineering 3.0 marks a paradigm shift in software development in which AI coding agents are no longer just assistive tools but active contributors. While prior empirical studies have examined productivity gains and acceptance patterns in AI-assisted development, the challenges associated with integrating agent-generated contributions remain less understood. In particular, **merge conflicts**, a fundamental aspect of collaborative software development, have received limited attention in this emerging context. In this paper, we present AGENTICFLICT, a large-scale dataset of textual merge conflicts in AI coding agent pull requests (Agentic PRs). The dataset comprises 142K+ Agentic PRs collected from 59K+ repositories, of which 107K+ are successfully processed through deterministic merge simulation. Our pipeline identifies 29K+ PRs exhibiting merge conflicts, yielding a conflict rate of 27.67%, and extracts 336K+ fine-grained conflict regions across these instances. Our preliminary exploratory analysis indicates that merge conflicts are both frequent and often substantial in AI-generated contributions, with noticeable variation across agents, emphasizing the need to better understand and manage integration challenges in AI-assisted software development.

The dataset supports analysis at multiple levels of granularity:

- Pull request level (conflict occurrence and summary statistics)
- File level (conflicting files and conflict types)
- Region level (fine-grained conflict structure with line boundaries)
- Repository level (project metadata)

All entities are linked using a canonical identifier `pr_key` (derived from `repo_full_name` and `pr_number`). Each record also includes a `run_id` to ensure traceability across dataset extraction runs.

A.2 Artifact Availability

- **Dataset (Zenodo):** <https://doi.org/10.5281/zenodo.19396916>
- **Preprint (arXiv):** <https://arxiv.org/abs/2604.03551>
- **Code and Documentation:** <https://github.com/unlv-evol/AgenticFlict>
- **API Documentation:** <https://agenticflict.readthedocs.io>

A.3 RAW vs. CLEAN Dataset

The replication package provides two versions of the dataset:

RAW Dataset. The RAW dataset contains the full output of the extraction pipeline, including:

- Pipeline metadata (`run_id`, timestamps, versioning)
- Simulation parameters and intermediate states (`simulation_mode`, `simulation_anchor_time`)

Authors' Contact Information: Daniel Ogenrwot, University of Nevada Las Vegas, Las Vegas, USA, ogenrwot@unlv.nevada.edu; John Businge, University of Nevada Las Vegas, Las Vegas, USA, john.businge@unlv.edu.

- Debug and error information (`status_code`, `error_class`, `error_message_trunc`)
- Optional raw conflict content (configurable at extraction time)

This version supports reproducibility, debugging, and re-execution of the pipeline.

CLEAN Dataset. The CLEAN dataset is a processed, analysis-ready version derived from the RAW dataset. It:

- Removes pipeline-specific metadata and debug fields
- Excludes raw source code content to comply with GitHub Terms of Service
- Retains only analysis-relevant attributes and join keys
- Ensures consistency across tables via `pr_key`

A detailed field-level mapping is provided in `raw_to_clean_manifest.csv` (included in the Zenodo archive). Unless otherwise stated, all results reported in the paper are based on the CLEAN dataset.

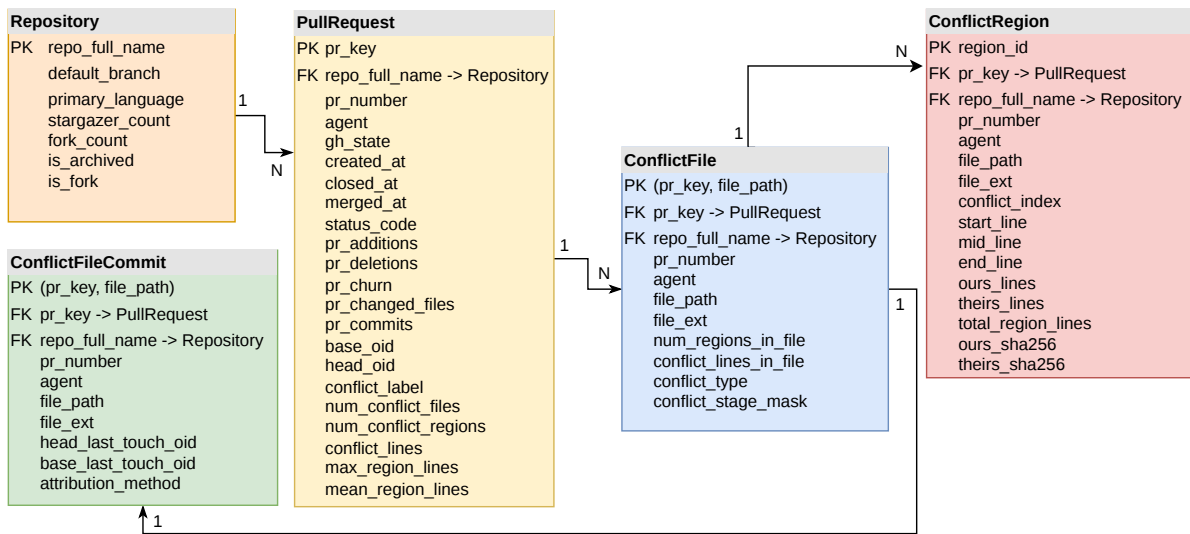


Fig. 1. Relational schema of the AGENTICFLICT CLEAN dataset. Tables are linked through the canonical `pr_key` identifier.

A.4 Dataset Files

The CLEAN dataset consists of the following tables:

- **`agenticflct_pr_clean.csv`** — Pull request-level metadata and conflict summary statistics. One row per simulated PR.
- **`agenticflct_conflict_files_clean.csv`** — File-level conflict information for each conflicting PR. One row per conflicting file.
- **`agenticflct_regions_clean.csv`** — Fine-grained conflict regions with line-level boundaries and cryptographic hashes. One row per conflict region.
- **`agenticflct_conflict_file_commits_clean.csv`** — Attribution of conflicting files to the most recent commits on each branch.
- **`agenticflct_repo_clean.csv`** — Repository-level metadata for all repositories in the dataset.

The RAW dataset contains the same tables with a `_raw` suffix, plus an additional table:

- **agenticflict_runlog_raw.csv** — Execution logs and failure tracking for every pipeline run.

A.5 Dataset Fields and Descriptions

Repository Table (*agenticflict_repo.csv*)

| Field | Description |
|-------------------------------|---|
| <code>repo_full_name</code> | Repository identifier (<code>owner/repo</code>) |
| <code>default_branch</code> | Default branch name |
| <code>primary_language</code> | Primary programming language |
| <code>stargazer_count</code> | Number of GitHub stars at extraction time |
| <code>fork_count</code> | Number of forks at extraction time |
| <code>is_archived</code> | Whether the repository is archived |
| <code>is_fork</code> | Whether the repository is itself a fork |

Pull Request Table (*agenticflict_pr.csv*)

| Field | Description |
|-----------------------------------|---|
| <code>pr_key</code> | Canonical identifier (<code>repo_full_name#pr_number</code>) |
| <code>repo_full_name</code> | Repository identifier (<code>owner/repo</code>) |
| <code>pr_number</code> | Pull request number within the repository |
| <code>agent</code> | AI coding agent that authored the PR |
| <code>gh_state</code> | GitHub PR state (<code>OPEN</code> or <code>CLOSED</code>) |
| <code>created_at</code> | PR creation timestamp (ISO 8601) |
| <code>closed_at</code> | PR closure timestamp (ISO 8601) |
| <code>merged_at</code> | PR merge timestamp; <code>NULL</code> for unmerged PRs |
| <code>pr_additions</code> | Number of lines added |
| <code>pr_deletions</code> | Number of lines deleted |
| <code>pr_changed_files</code> | Number of files modified |
| <code>pr_commits</code> | Number of commits in the PR |
| <code>pr_churn</code> | Total code churn (<code>additions + deletions</code>) |
| <code>base_oid</code> | Base branch commit SHA (from GitHub) |
| <code>head_oid</code> | Head branch commit SHA (from GitHub) |
| <code>simulation_mode</code> | Anchor strategy: <code>open_current</code> or <code>closed_base_at_close</code> |
| <code>conflict_label</code> | Boolean: <code>True</code> if merge simulation produced a textual conflict |
| <code>num_conflict_files</code> | Number of files containing conflict markers |
| <code>num_conflict_regions</code> | Total conflict regions across all files |
| <code>conflict_lines</code> | Total lines enclosed by conflict markers |
| <code>max_region_lines</code> | Lines in the largest single conflict region |
| <code>mean_region_lines</code> | Mean lines per conflict region |

Conflict Region Table (agenticflict_regions.csv)

| Field | Description |
|--------------------|---|
| region_id | Unique identifier for the conflict region |
| pr_key | Foreign key to the PR table |
| repo_full_name | Repository identifier |
| pr_number | Pull request number |
| agent | AI coding agent |
| file_path | Relative path of the conflicting file |
| file_ext | File extension |
| conflict_index | 0-based index of the region within the file |
| start_line | Line number of the <<<< marker |
| mid_line | Line number of the ===== separator |
| end_line | Line number of the >>>> marker |
| ours_lines | Number of lines on the base (ours) side |
| theirs_lines | Number of lines on the incoming (theirs) side |
| total_region_lines | Total lines in the conflict region |
| ours_sha256 | SHA-256 hash of the ours conflict block |
| theirs_sha256 | SHA-256 hash of the theirs conflict block |

Conflict File Table (agenticflict_conflict_files.csv)

| Field | Description |
|------------------------|--|
| pr_key | Foreign key to the PR table |
| repo_full_name | Repository identifier |
| pr_number | Pull request number |
| agent | AI coding agent |
| file_path | Relative path of the conflicting file |
| file_ext | File extension |
| num_regions_in_file | Number of conflict regions in this file |
| conflict_lines_in_file | Total lines enclosed by conflict markers in this file |
| conflict_type | Git 3-way merge classification (e.g., both_modified, both_added) |
| conflict_stage_mask | Encoded unmerged index stage information |

Conflict File Commit Table (*agenticflict_conflict_file_commits.csv*)

| Field | Description |
|----------------------------------|---|
| <code>pr_key</code> | Foreign key to the PR table |
| <code>repo_full_name</code> | Repository identifier |
| <code>pr_number</code> | Pull request number |
| <code>agent</code> | AI coding agent |
| <code>file_path</code> | Relative path of the conflicting file |
| <code>file_ext</code> | File extension |
| <code>head_last_touch_oid</code> | SHA of the last commit to touch this file on the head branch |
| <code>base_last_touch_oid</code> | SHA of the last commit to touch this file on the base branch |
| <code>attribution_method</code> | Method used for attribution (<code>git_log_last_touch</code>) |

B Extraction Pipeline and Reproducibility

B.1 Pipeline Overview

The extraction pipeline processes each Agentic PR through seven sequential steps:

- (1) **Source loading.** Agentic PRs are loaded from the hao-li/AIDev HuggingFace dataset, which catalogues AI agent-authored pull requests.
- (2) **Metadata retrieval.** Repository and PR metadata are fetched via the GitHub GraphQL API, including commit OIDs, PR state, and repository attributes.
- (3) **Filtering.** Merged PRs are excluded from the dataset; only OPEN and CLOSED-unmerged PRs are retained.
- (4) **Repository caching.** Repositories are cloned locally using blobless shallow clones (`-filter=blob:none`) and cached on disk.
- (5) **Simulation anchor selection.** For OPEN PRs, the current `base_oid` and `head_oid` are used. For CLOSED-unmerged PRs, the base branch state at `closed_at` is resolved via `git rev-list`.
- (6) **Merge simulation.** A local `git merge -no-commit -no-ff` is executed using the anchored commit OIDs.
- (7) **Conflict extraction.** Conflict markers (`<<<</=====/>>>>`) are parsed to extract region-level metrics and SHA-256 hashes.

B.2 Reproducing a Single Label

Each PR row stores `base_oid` and `head_oid`. Any conflict label can be independently verified from public GitHub history:

```

1 # Clone the repository
2 git clone https://github.com/{owner}/{repo} /tmp/verify-repo
3 cd /tmp/verify-repo
4
5 # Fetch the required commits and run the merge simulation
6 git fetch origin {head_oid}
7 git checkout {base_oid}
8 git merge --no-commit --no-ff {head_oid}
9
10 # Non-zero exit code and conflict markers confirm a conflicting PR

```

Listing 1. Verifying a single conflict label

C Loading and Using the Dataset

C.1 Requirements

- Python 3.9 or higher
- pandas (pip install pandas)

C.2 Loading the CLEAN Dataset

```

1 from pathlib import Path
2 import pandas as pd
3
4 DATA_DIR = Path("data/clean")
5
6 pr          = pd.read_csv(DATA_DIR / "agenticflict_pr_clean.csv")
7 regions    = pd.read_csv(DATA_DIR / "agenticflict_regions_clean.csv")
8 cfiles     = pd.read_csv(DATA_DIR / "agenticflict_conflict_files_clean.csv")
9 commits    = pd.read_csv(DATA_DIR / "agenticflict_conflict_file_commits_clean.csv"
10                      )
11 repo       = pd.read_csv(DATA_DIR / "agenticflict_repo_clean.csv")
12 print(f"PRs: {len(pr):,} | Conflicting: {pr['conflict_label'].sum():,}")

```

Listing 2. Loading the CLEAN dataset

C.3 Loading the RAW Dataset

```

1 from pathlib import Path
2 import pandas as pd
3
4 DATA_DIR = Path("data/raw")
5
6 pr          = pd.read_csv(DATA_DIR / "agenticflict_pr_raw.csv")
7 regions    = pd.read_csv(DATA_DIR / "agenticflict_regions_raw.csv")
8 cfiles     = pd.read_csv(DATA_DIR / "agenticflict_conflict_files_raw.csv")
9 commits    = pd.read_csv(DATA_DIR / "agenticflict_conflict_file_commits_raw.csv")
10 repo       = pd.read_csv(DATA_DIR / "agenticflict_repo_raw.csv")
11 runlog     = pd.read_csv(DATA_DIR / "agenticflict_runlog_raw.csv")
12
13 print(f"PRs: {len(pr):,} | Run log entries: {len(runlog):,}")

```

Listing 3. Loading the RAW dataset

C.4 Example Analyses

```

1 agent_stats = (
2     pr.groupby("agent")["conflict_label"]

```

```

3     .agg(total="count", conflicts="sum")
4     .assign(conflict_rate=lambda d: d["conflicts"] / d["total"])
5     .sort_values("conflict_rate", ascending=False)
6 )
7 print(agent_stats)

```

Listing 4. Conflict rate by AI agent

```

1 merged = regions.merge(
2     pr[["pr_key", "agent", "repo_full_name"]], on="pr_key"
3 )
4 print(merged.groupby("agent")["total_region_lines"].describe())

```

Listing 5. Joining PR metadata with conflict regions

C.5 Notes

- The CLEAN dataset is recommended for analysis and modeling.
- The RAW dataset includes additional pipeline metadata and execution logs useful for debugging and reproducibility audits.
- All tables can be joined using `pr_key` as the primary identifier.
- Repository metadata reflects values at extraction time and may differ from current repository state.

D Dataset Usage

AGENTICFLICT is designed to support, but is not limited to, the following research directions:

- **Conflict prediction.** Train classifiers that predict whether an incoming agent PR will conflict, using PR size, repository activity, file overlap, and agent identity as features.
- **Agent behavior comparison.** Compare conflict profiles (rate, severity, type) across the five included AI agents to identify behavioral patterns that elevate integration risk.
- **Automated conflict resolution.** Use paired conflict-region data (including SHA-256 hashes of both sides of each conflict) as inputs to resolution models.
- **Software repository mining.** Study how repository characteristics (language, stars, fork count) moderate conflict likelihood in AI-authored contributions.
- **Longitudinal analysis.** Use `created_at` and `closed_at` timestamps to analyze how agent deployment patterns and conflict rates have changed over time.

E Limitations

- **Textual conflicts only.** The dataset captures syntactic merge conflicts as detected by Git's 3-way merge. Semantic conflicts (syntactically correct but behaviorally incompatible changes) are not detected.
- **Excluded PRs.** PRs are excluded when repositories are deleted, made private, or inaccessible at extraction time, introducing potential survivorship bias toward more active repositories.
- **No human-authored baseline.** The current version focuses exclusively on AI agent PRs. A matched human-authored baseline will be added in a future release to enable direct comparison.
- **Snapshot in time.** Conflict labels for OPEN PRs reflect repository state at extraction time. Re-running the pipeline at a later date may yield different results as base branches advance.

8 ■ Daniel Ogenrwot and John Businge

- **Agent coverage.** Coverage is limited to agents represented in the hao-li/AIDev source dataset. Agents with few PRs may exhibit higher variance in per-agent metrics.