

# WHY DO YOU FAIL ME MR. BOT?

2025/10/28

Amaan Ahmed, Asim Waheed, Youssef Souati

CS846 – Project Presentation



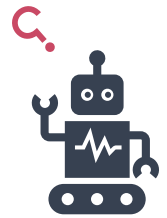
## Analyzing Pull-Request Failures in the AIDev Dataset

AI Coding Agents now author hundreds of thousands of PRs

But many PRs **fail to be merged** or are **reverted**

**Our goal:**

Understanding *why*



PAGE 2



## Why study PR failures?

PR acceptance = natural test of **trust** in AI teammates

Prior work<sup>[1]</sup> showed agent PRs accepted **far less** than humans

Recent data suggests **trend reversed**

**Understanding failures** → **Improving human-AI workflows**

## The AIDev Dataset

- 900,000+ agent-authored PRs across 100,000+ repos
- Metadata on PRs, reviews, comments and timelines
- Rich comparisons between agents and humans
- Useful to study PR failure across authors and contexts

[1] Wyrick, M., Ghit, R., Haller, T., & Miller, C. (2021). Bots Don't Mind Waiting, Do They? Comparing the Interaction With Automatically and Manually Created Pull Requests. 2021 IEEE/ACM Third International Workshop on Bots in Software Engineering (BotSE), 6–10. <https://doi.org/10.1146/robot.2021.000009>

## Research Questions

**RQ1:** Are agent-authored PRs more likely to succeed than human-authored PRs?

**RQ2:** What factors lead to **rejection** of agentic PRs?

**RQ3:** What practices make agentic PRs more **successful**?



PAGE 5

## METHODOLOGY

What even is a “failed” PR?

PAGE 6

## Dataset Details

- Datasets:
  - Curated AIDev Dataset: 33,596 PRs.
- Definition of **failed** PRs.
  - PR state is closed and *merged\_at* is NULL.
  - PR state is open but no activity for 180 days.
- 7,646 failed PRs comprising 7,270 closed-but-unmerged and 376 stale.

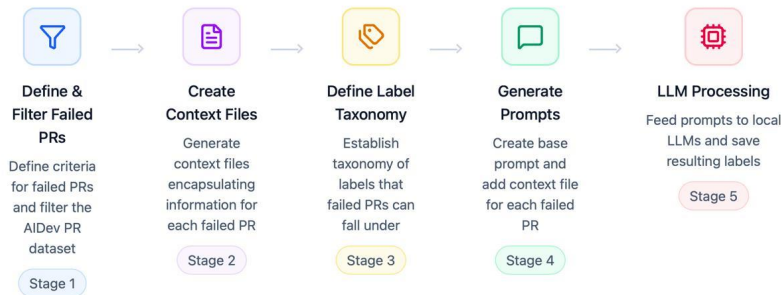
PAGE 7

## Research Question 1

- Extract 6,618 human PRs from curated dataset
- Apply same failure rule → 1,204 failed PRs.
- Compare task-type distributions and failure rates by task type.

PAGE 8

## Research Question 2



## RQ2: Stage 2

- PR context files include:
  - Repository metadata
  - Diff summary
  - Event timeline
  - CI/workflow signals
  - Review text
  - Commit messages
  - Final PR outcome

```

Pin #2756921963
Repository: unknown (Lang: Python)
Outcomes: closed_unmerged
Signals: CI_fail=False, Changes_Req=False

Title: Add LLM Integration Tests

Description:
# @ Review Summary
==@reviewsummary==
- Enhance the testing framework by integrating LLM providers.
- Add tests for asynchronous and streaming LLM interactions.
- Improve error handling and logging for LLM-related operations.

==@changes==
- @Configurations: Introduced environment variables for LLM API keys and providers.
- @Enhancements: Improved asynchronous handling of LLM requests.
- @Tests: Expanded testing to include both synchronous and asynchronous LLM interactions.
- @Miscellaneous: Added logging for LLM API calls and responses.

==@misc==
- AI21
- AI21

==@discussions==
This pull request has been automatically marked as stale because it has not had any recent activity. It will be closed if no updates are made within 7 days, this PR will be automatically closed.

==@devin AI Engineer==
The PR also:
- Adds necessary test dependencies to the testing workflow with required packages.
- Adds debug prints for API key and secret values.
- Adds a link to Devin run: https://app.devin.ai

I'll be helping with this pull request! Here's what you should know:
- I will automatically:
  - Add comments on this PR. Add "aside" to your comment to have me ignore it.
  - Look at CI failures and help fix them.

@ Control Options:
- [ ] Disable automatic comment and CI monitoring

==@details==
File: github/workflows/python-testing.yml
@ -24,6 +24,12 @ jobs:
  test:
    runs-on: ubuntu-latest

```

## RQ2: Stage 3 (Failure Taxonomy)

Group	Description	Category
TECHNICAL (T)	Code, tests, style, architecture	T1: CI or Test Failure T2: Code Quality or Correctness
FTT_OR_VALUE (F)	Feature/requirement mismatch, not needed, wrong direction	F1: Wrong feature, not needed, or misalignment
PROCESS_OR_OWNERSHIP (P)	Time, priority, unresponsiveness, social/review issues, policy	P1: Low priority P2: Author Unresponsive P3: Review/social coordination conflict P4: Policy/compliance/bureaucracy
REDUNDANCY_OR_OBSOLETE (R)	Duplicate feature or superseded work	R1: Duplicate feature R2: Obsolete FR
INVALID_OR_SPAM (S)	Spam, noise, non-meaningful PRs	S1: Invalid/Spam/Nonsensical
UNKNOWN (U)	Not enough information	U1: Unknown

Table 1: Failure taxonomy used for classifying rejected PRs

```
{
  "id": "T1_CI_OR_TEST_FAILURE",
  "group": "TECHNICAL",
  "definition": "The PR is not merged mainly
    because tests, builds, or CI checks fail due
    to issues introduced or exposed by the PR, and
    the author never fixes them.",
  "typical_signals": [
    "ci_failed = true",
    "Timeline events like 'ci.failure', '
    workflow_run.completed: failure'",
    "Review comments: 'tests are failing', 'please
    fix ci', 'build is broken'",
    "No later commit that fixes the failures
    before close"
  ]
}
```

## RQ2: Stage 4 and 5

- Base prompt + context file = prompt for each failed PR
- Output: (label, justification, confidence)
- Output JSON-enforced
- [LLaMa-3.2-3b](#) and [Qwen2.5](#) used

You are an assistant that specializes in how GitHub prompt requests failed.

A "failed" prompt is one that has been used that was CLOSETED without being used, or has remained OPEN and STALE for a long time.

Your goal is to read the prompt request content and decide the SINGLE most likely primary reason why it did not work, according to the taxonomy below.

If there is *generally* not enough information in the request to infer a reason, you MUST use the UNKNOWN label (U1,UNKNOWN).

Do not guess when the evidence is missing or extremely weak.

Follow the taxonomy <CHOOSE ONE>:

TASK:

1. Read the PR content carefully.
2. Decide the SINGLE most likely primary reason label from the taxonomy above.
3. Write a short reason summary (2-3 sentences) explaining why you chose that label.
4. Estimate a numeric confidence between 0 and 1 for how sure you are.

IMPORTANT:

- Choose **exactly ONE** UNKNOWN label.
- Do NOT insert details that are not supported by the PR.
- If the signals are mixed, choose the reason that best explains why the PR ultimately failed.
- If there is not enough information, use U1,UNKNOWN.

OUTPUT FORMAT (JSON ONLY):

You MUST output exactly ONE JSON object with this schema:

```
{
  "reason": "U1,UNKNOWN",
  "reason_label": "U1,UNKNOWN",
  "reason_summary": "412 sentence explanations",
  "confidence": "0.1 between 0 and 1"
}
```

All string values (reason label and reason summary) should be enclosed in double quotes " ".

pr id must match the pr id given below.

- pr id MUST include any text before or after the JSON.
- DO NOT include any extra text.
- DO NOT include any text before or after the JSON.

Now read the PR content below and produce the JSON output.

PR Content: <pr content>

## Research Question 3

- Combine the following data sources:
  - Failure reasons from RQ2
  - PR-level metadata (PR size, number of reviewers, etc.)
  - Repo-level metadata (Age, stars, etc.)
- Analyze failure reasons grouped by each variable
  - Report variables with **significant differences** for different values

PAGE 13

# RESULTS

So, why DO you fail me, Mr. Bot?

PAGE 14

## Recap

- Questions to answer:
  - Do agents **fail more** than humans?
  - **Why** do agent PRs fail?
  - What **practices** and **contexts** make agents **more likely to succeed**?
- All numbers here from curated dataset:
  - ~33K agent PRs.
  - ~6.6K human PRs.
  - Popular, actively maintained repos.

PAGE 15

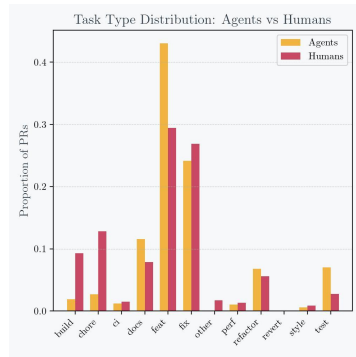
## RQ1 – Do Agents Fail More Often Than Humans?

- Using curated repos:
  - **Agents**: 22.79% failure rate
  - **Humans**: 19.70% failure rate
- Using full AIDev dataset agents perform better (**8.18% failure**)
  - Due to less popular repos and less scrutiny during review.

In **serious, high-activity** projects, agents **fail more often than** humans.

PAGE 16

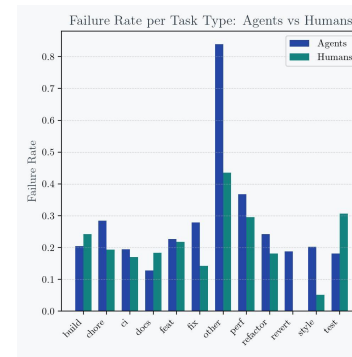
## RQ1 - What Tasks Do Agents vs Humans Work On?



- Agents handle generative or patch-style work:
  - New feature
  - Documentation
  - Tests
- Humans handle infrastructure & maintenance
  - Requires more general context

PAGE 17

## RQ1 - Where Do Agents Fail vs Humans?



- Agents
  - Strong on structured, repetitive tasks, e.g., docs and tests
  - Struggle on context-heavy maintenance work, especially code style
- Humans
  - Very strong on style and atypical PRs ('other' category)
  - Struggle with test-writing and docs

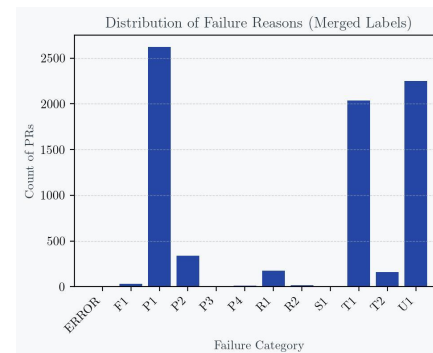
PAGE 18

## RQ2: Evaluating model agreement

- Models agree for only ~45% labels.
- Upon analyzing:
  - Qwen2.5 tends to classify ambiguous failures as P1 (low priority).
  - Structured disagreement → systemic bias.
- Fix:
  - When both models label a failure as P1 → true label is P1.
  - If one model labels failure as P1 → pick other model's label.

PAGE 19

## RQ2 - Classifying Why Agent PRs Fail



- Most Agent PRs fail due to low priority.
- Technical failures result from CI failure.
- Social conflicts and spam are negligible.

PAGE 20

## RQ3 – How Task Type, Size, and Repo Context Affect Success

- Task type vs failure reason:
  - **Low-impact** tasks tend to **fail silently**.
  - **Code-changing** tasks fail mainly at the CI gate.
- PR Size and Architectural Footprint:
  - Agents perform better when PRs are **narrow** in scope and **easy** to review.
  - **Large, single-commit diffs** across many files **more likely to fail**.
- Repository Characteristics:
  - **Agents struggle** in **large, high-visibility** projects with **strict review** and **lots of competing PRs**.

PAGE 21

## WHAT NOW?

PRESENTATION TITLE

PAGE 22

## Future Work

- Using better LLMs for classification
- Repeating study across richer datasets
- Creating manually-labelled dataset for failure reasons
- Improving taxonomy to minimize unknown labels

PRESENTATION TITLE

PAGE 23

## Takeaways and Recommendations

- **Agents** perform best when:
  - Project scope and priority is **clearly defined**
  - **Repetitive tasks** such as boilerplate code, documentation, writing tests
  - PRs are **small, narrowly scoped**, and include **fewer changes**
- **Humans** perform best when:
  - Working on **atypical problems**
  - Tasks requiring **lots of context** (performance enhancements, CI configuration)
- Most PR failures result from **low priority** or **CI failure**

PAGE 24