# Automated GitHub Repository Quality Evaluation: A Metrics-Based Approach

Dmytro Polishchuk

# Github overview

Ex. https://github.com/ansible/ansible

<mark>1. Commits</mark>
Commits represent individual changes to a repository's codebase. Each commit records a snapshot of the project's state at a specific point in time.
<mark>2. Issues</mark>
Issues represent tasks, bugs, or feature requests within a repository. They are often used to track the progress of software development and maintain a record of discussions around various problems or enhancements.
https://github.com/ansible/ansible/issues/84909
<mark>3. Pull Requests (PRs)</mark>
Pull requests are used to propose changes to a repository. They are essential for collaborative workflows, as they allow other contributors to review code changes before merging them into the main codebase.
https://github.com/ansible-community/antsibull-nox/pull/26
<mark>4. Actions</mark>
GitHub Actions is a CI/CD tool that allows developers to automate workflows, such as running tests, building applications, and deploying software when certain events occur in a repository.
https://github.com/ansible-community/antsibull-nox/pull/26/checks
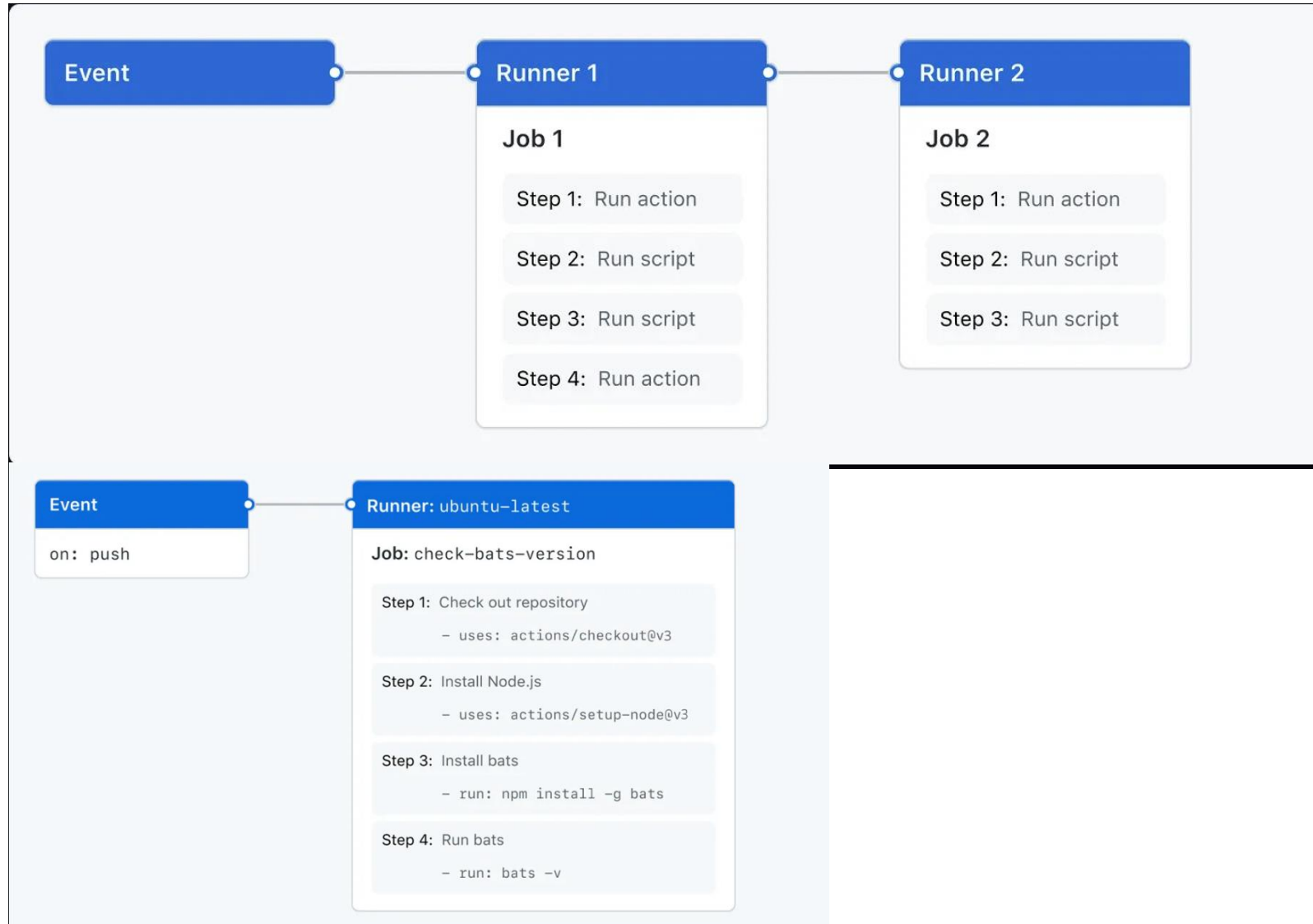<mark>5. Timeline</mark>
The GitHub timeline is a chronological sequence of events that provides an overview of all activities related to an issue or pull request. It captures all updates, comments, status changes, and relevant interactions.
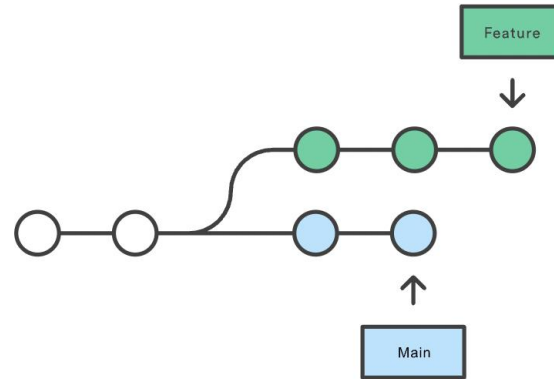https://github.com/ansible-community/antsibull-nox/pull/26/commits/8b7f4e2dd368c7aa471bffd2ca6a3f7c5f918408
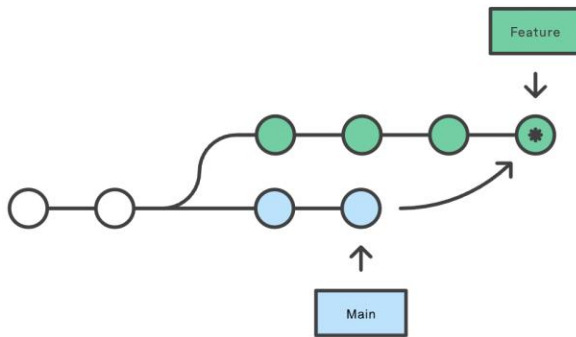
# Github workflows:
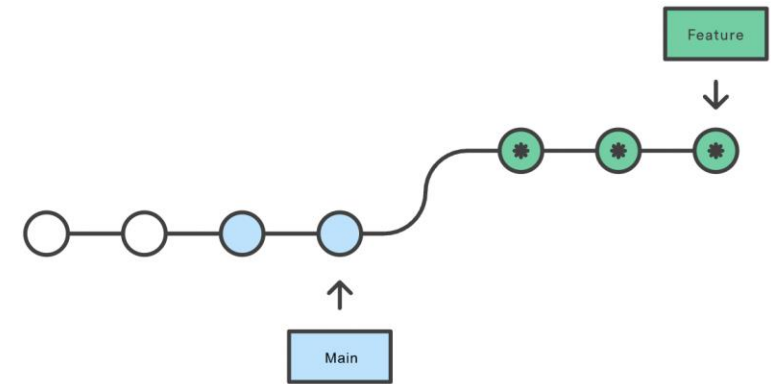
Possible scenarios:

1. Merge main into feature

2. Rebase feature onto main

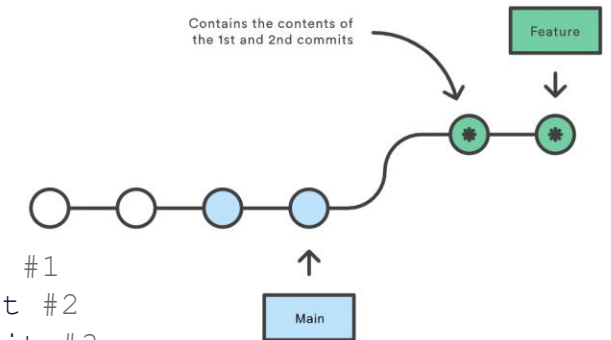git checkout feature
git rebase main



git checkout feature
git merge main



```
pick 33d5b7a Message for commit #1
pick 9480b3d Message for commit #2
pick 5c67e61 Message for commit #3
```

Interactive rebase
git rebase -i main

Contains the contents of
the 1st and 2nd commits



```
pick 33d5b7a Message for commit #1
fixup 9480b3d Message for commit #2
squash 5c67e61 Message for commit #3
```

# Reset/Checkout/Revert



Before Resetting

Hotfix / HEAD

Main

Before Checking Out

Hotfix

Main / HEAD

Before Reverting

Hotfix

Main

Commit to be reverted

After Resetting

Hotfix / HEAD

Main

After Checking Out

Hotfix / HEAD

Main

After Reverting

Removes changes from the reverted commit

Feature

Main

# Git Hooks

Maintaining a hook using a symlink to version-controlled script

Hook

Symlink

Hooks

User Actions

.git/hooks

Git Repo

Pre-Commit
Notification

Post-Commit
Notification

1
Stage
Changes

2
Commit
Changes

3
Enter Commit
Message

4
All Done!

# GitHub WebHooks

GitHub

git push

Webhook
event

git pull&&
build.sh

developer

server

Best practices:

Reviewing a feature with a pull request

If you use pull requests as part of your code review process, you need ==to avoid using git rebase after creating the pull request==. As soon as you make the pull request, other developers will be looking at your commits, which means that it's a *public* branch. Re-writing its history will make it impossible for Git and your teammates to track any follow-up commits added to the feature. Any changes from other developers need to be incorporated with git merge instead of git rebase. For this reason, it's usually a good idea to clean up your code with an interactive rebase *before* submitting your pull request.

Integrating an approved feature

After a feature has been approved by your team, you have the option of rebasing the feature onto the tip of the main branch before using git merge to integrate the feature into the main code base.

This is a similar situation to incorporating upstream changes into a feature branch, but since you're not allowed to re-write commits in the main branch, you have to eventually use git merge to integrate the feature. ==However, by performing a rebase before the merge, you're assured that the merge will be fast-forwarded==, resulting in a perfectly linear history. This also gives you the chance to squash any follow-up commits added during a pull request.

## 2. Related Work and Background

## 2.1 Software Quality Models

Assessing software quality is a central concern in software engineering research, with various models developed to quantify and measure different aspects of quality. Traditional approaches, such as the ISO/IEC 25000:2014 (https://www.iso.org/standard/64764.html) series (derived from ISO/IEC 25002:2024 https://www.iso.org/standard/78175.html), provide structured quality evaluation frameworks focused on attributes like maintainability, reliability, and usability. These models define software quality from a product-oriented perspective but often overlook the dynamic aspects of software development, such as the role of development workflows, community activity, and the structure of version control histories. Several models extend these traditional frameworks to better suit the open-source software (OSS) landscape. The "Measuring the Quality of Open Source Software Ecosystems Using QuESo" (https://www.researchgate.net/publication/300335654_Measuring_the_Quality_of_Open_Sou rce_Software_Ecosystems_Using_QuESo) model was introduced to address sustainability and maturity in open-source projects by considering both technical and community-related factors. The model highlights the importance of process stability, development continuity, and project governance in assessing long-term software quality. Similarly, QuESo expands the scope by incorporating ecosystem-level quality metrics, including project interdependencies, contributor engagement, and software reuse. These ecosystem oriented models provide valuable insights into software quality at a higher level but do not focus on how specific development structures such as commit history topologies affect software maintainability and defect rates. The role of software ecosystems in ensuring the sustainability of OSS projects has been widely studied. The " QuESo" introduces a comprehensive framework for evaluating OSS ecosystems based on multiple quality dimensions, such as activity levels, stability, governance, and modularity. The study highlights that a well-structured and active software ecosystem contributes to the long live stability and reliability of OSS projects, as continuous contributions and frequent interactions between developers help maintain a project's health. Another relevant study, "Quality Evaluation Meta-Model for OSS – A Multi-Method Validation Study" https://www.researchgate.net/publication/378003488_Quality_evaluation_meta- model_for_open-source_software_multi-method_validation_study presents a meta-model for software quality assessment, creating multiple prior approaches into a unified framework. The study emphasizes the importance of historical project data, including commit activity, developer collaboration, and repository structure, as indicators of software sustainability and robustness. However, while these studies explore high-level ecosystem characteristics and development patterns, they do not directly investigate the commit structure as a directed acyclic graph (DAG). The DAG structure in version control systems encodes essential information about how a project evolves over time, capturing branching strategies, merge practices, and code integration workflows. Understanding how these structural properties correlate with software quality can provide new insights into best practices for OSS development.

Data Set Collection

To study software evolution, bug-fixing patterns, and code maintenance, we need **real-world data**.

Without a dataset, research becomes **theoretical** and lacks empirical evidence (other researches are based on expert opinion in building evaluation model).
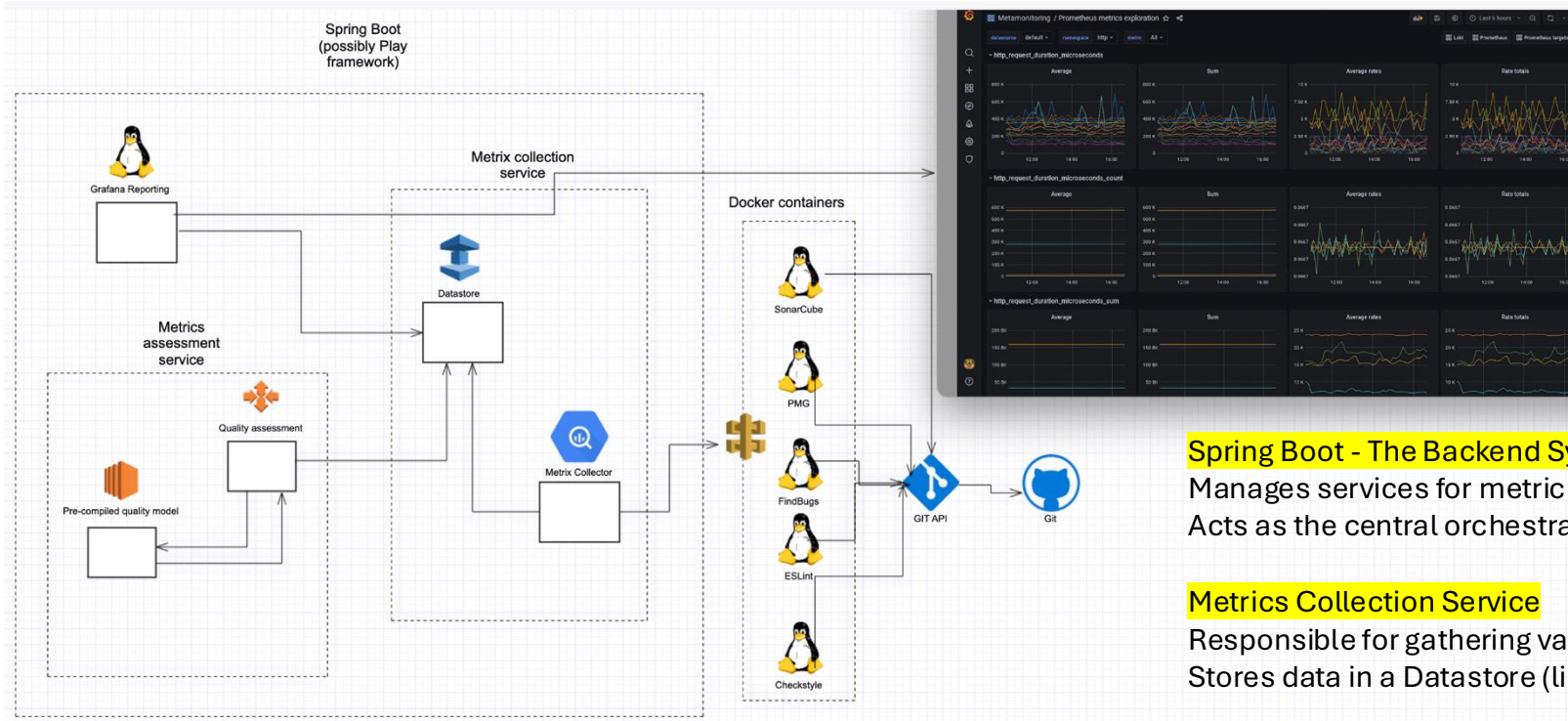
To ensure a comprehensive dataset, I selected:

- Top 5 GitHub projects based on popularity, activity, and long-term maintenance.
- Projects with at least 15+ years of history to capture long-term software evolution.
- A proportional number of pull requests (PRs) and issues, ensuring a balanced dataset with various software development patterns.
- Preference for repositories with high collaboration and structured development workflows.

Collecting data from GitHub is not straightforward due to API constraints:

- The GitHub API limits requests to 1,000 per hour per user, making large-scale data collection difficult.
- Workarounds such as caching responses and rate limit handling were implemented.
- Orphaned Commits & Forks
- Some commits do not belong to any branch, making them difficult to trace (orphaned commits).
- Implemented an enhanced SZZ algorithm to recover these orphaned commits.
- Data Integrity & Completeness
- Many PRs do not directly link to issues, requiring heuristics to infer bug-fixing commits.
- Actions & workflows in CI/CD pipelines add additional complexity in identifying changes triggered by bug fixes.

# Architectural design:



## Spring Boot - The Backend System
Manages services for metric collection and quality assessment.
Acts as the central orchestrator of the pipeline.

## Metrics Collection Service
Responsible for gathering various software metrics.
Stores data in a Datastore (likely a database or time-series storage).

Docker Containers - Static Code Analysis Tools
SonarCube, PMG, FindBugs, ESLint, Checkstyle run inside Docker containers.
These tools analyze the code and generate reports on code quality, security, and bugs.

## Metrics Assessment Service
Quality assessment component evaluates the collected data.
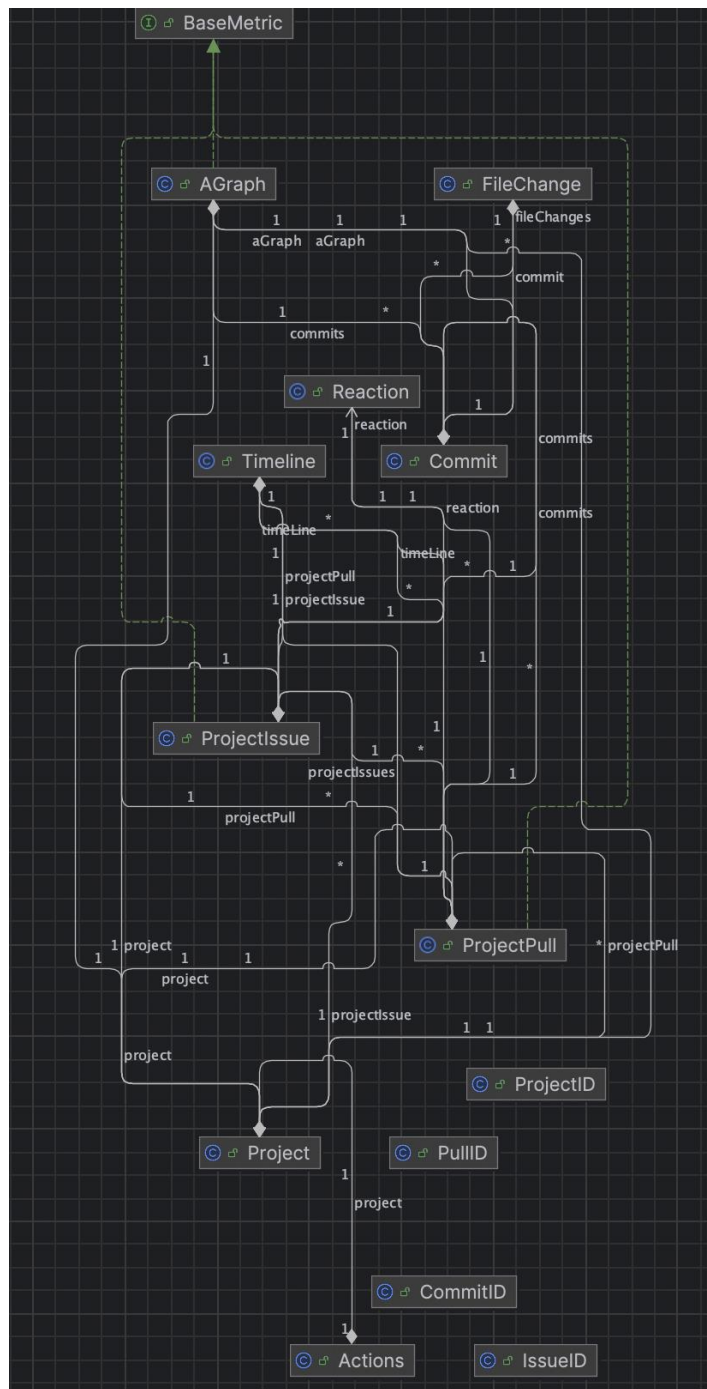Uses a pre-compiled quality model to assess software quality.

## Git API & GitHub Integration
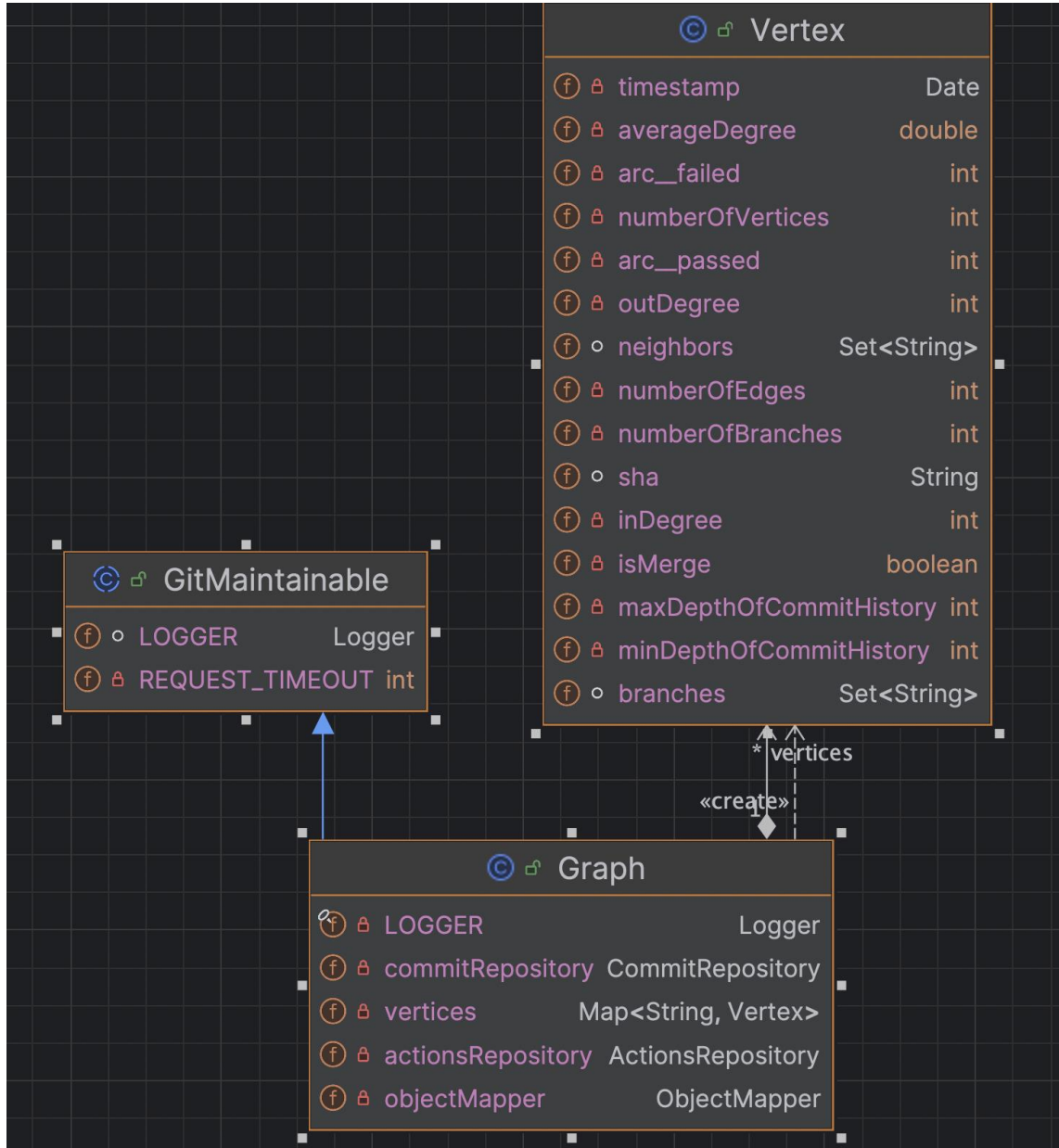Git API fetches repository data (e.g., commits, PRs, issues).

## Reporting
Presents visual insights on collected data.

# Proposed Schema

# ORM view

# DAG of commits (possibly demo)

**Vertex**
| | |
|---|---|
| timestamp | Date |
| averageDegree | double |
| arc_failed | int |
| numberOfVertices | int |
| arc_passed | int |
| outDegree | int |
| neighbors | Set<String> |
| numberOfEdges | int |
| numberOfBranches | int |
| sha | String |
| inDegree | int |
| isMerge | boolean |
| maxDepthOfCommitHistory | int |
| minDepthOfCommitHistory | int |
| branches | Set<String> |

**GitMaintainable**
| | |
|---|---|
| LOGGER | Logger |
| REQUEST_TIMEOUT | int |

**Graph**
| | |
|---|---|
| LOGGER | Logger |
| commitRepository | CommitRepository |
| vertices | Map<String, Vertex> |
| actionsRepository | ActionsRepository |
| objectMapper | ObjectMapper |

\* vertices

«create»

**Adding Commits (Vertices):**
Each commit (represented by SHA) is added as a vertex in the graph.
The Graph maintains a collection of Vertices (commits), each containing a unique SHA, neighbors (other commits that depend on it), branches (where it belongs), and degree metrics.
Building Commit Relationships (Edges):
When a commit has parent commits (e.g., for merges), the Graph will add edges between these commits.
The addEdge method establishes these relationships and updates the in-degree and out-degree properties of the vertices.

**Tracking Commits:**
The findParents method finds all parent commits for a given commit (sha).
The findChildren method finds all child commits.
The findCommonAncestors method finds common ancestors between multiple commits, helping identify shared history in the commit graph.

**Calculating Metrics:**
The Graph computes metrics related to the graph structure (e.g., number of commits, branches, edges).
Vertex tracks the number of neighbors (out-degree) and the number of incoming edges (in-degree), helping measure graph complexity.

**Commit Analysis:**
The system analyzes commit depth (max/min depth), branching strategies, and merge count, which are stored as properties in each Vertex.
The updateActionResult method links build results (passed/failed) with commits, updating the commit graph accordingly.

# Example of graph for ansible project

# Commits that fix the issue:



Ex. Find bug fixing commits
https://github.com/facebook/react/issues/10
https://github.com/facebook/react/pull/11
https://github.com/facebook/react/pull/11/commits

# SZZ Algorythm (Software Bug-Inducing Commit Identification)

The algorithm tracks bug-fixing commits (commits that fix an issue) and then traces backward through the commit history to identify the bug-introducing commits

# Handle orphaned commits:



Orphaned commits example:
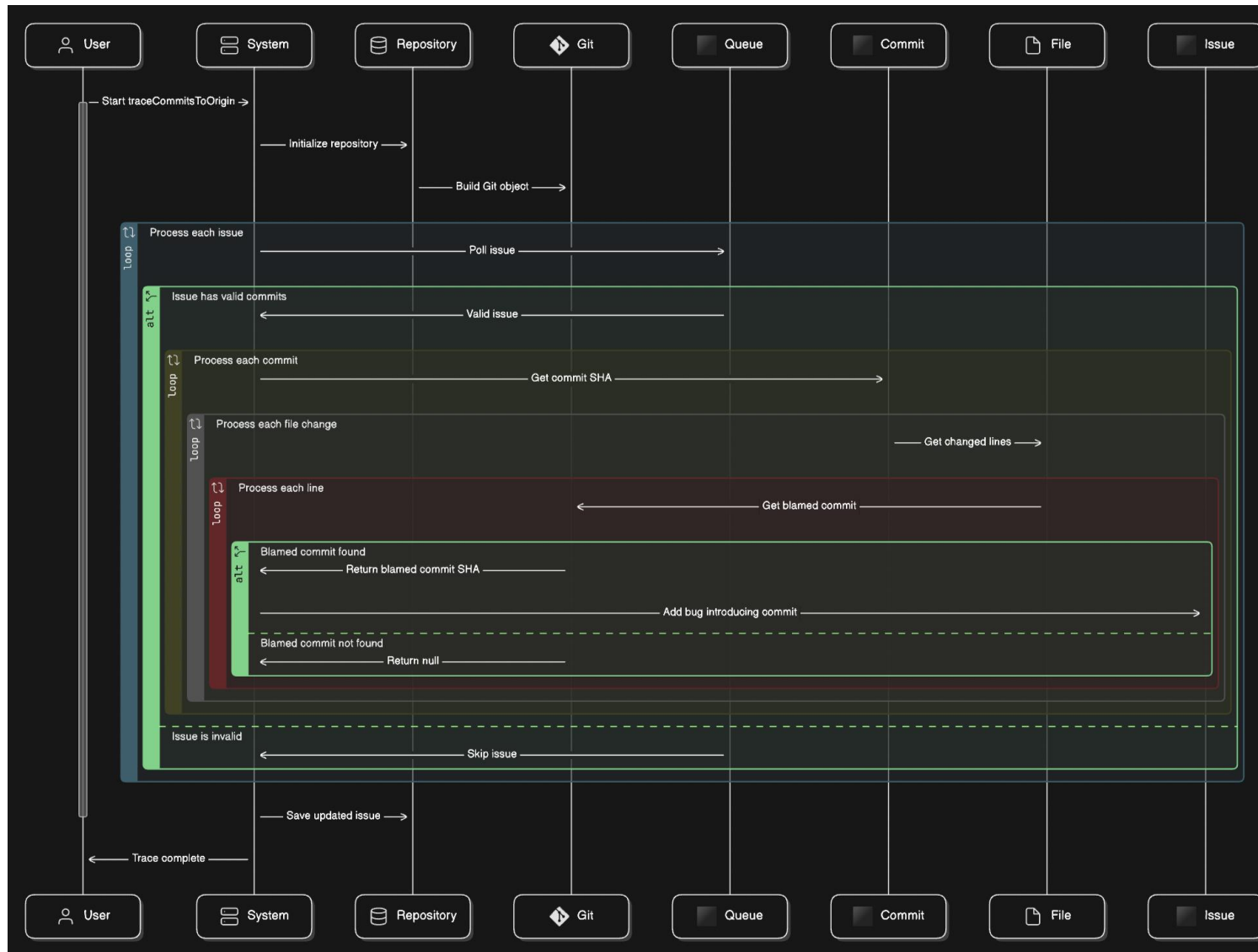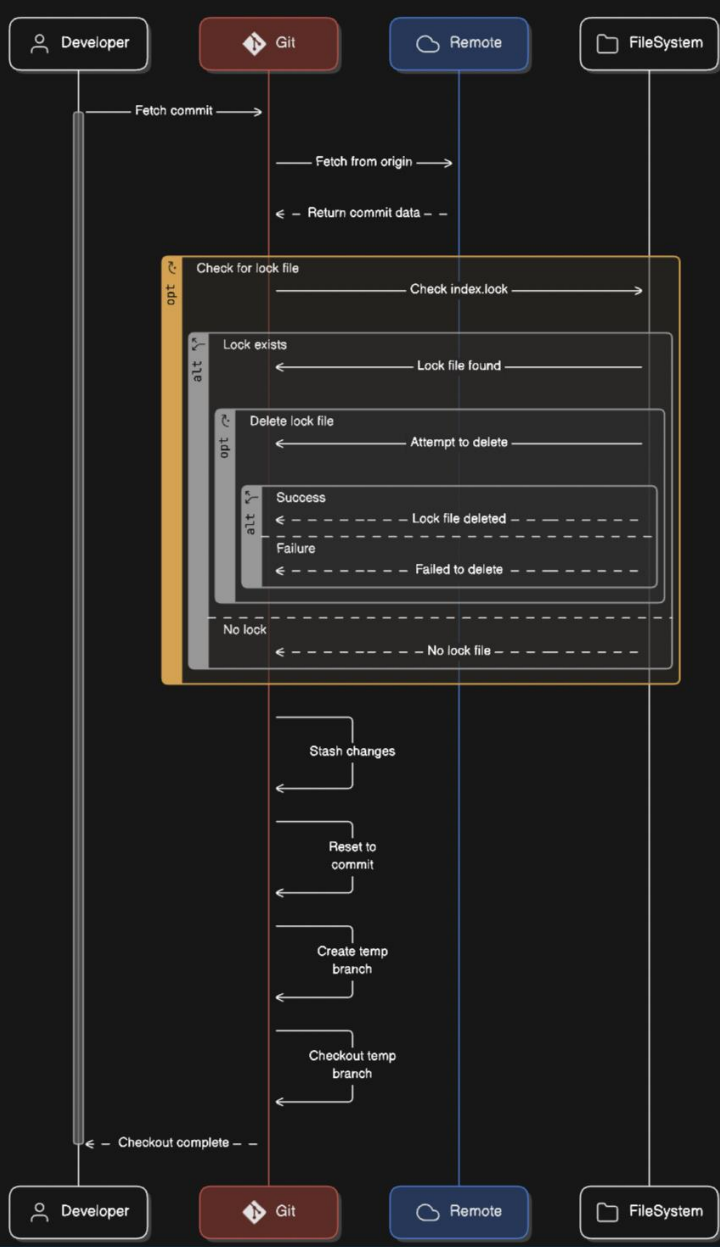https://github.com/google/guava/commit/0ceaed084bb7e52d8fa3b0760cb0ffe486918a00

# Working with the dataset

select concat(project_owner, '_', project_project_name) as project, date(created_at) as time, count(*) as issues_opened from project_issue group by project_owner, project_project_name, date(created_at) order by date(created_at) asc;

SELECT CONCAT(project_owner, '_', project_project_name) AS project, DATE(closed_at) AS time, COUNT(*) AS issues_closed FROM project_issue WHERE closed_at IS NOT NULL GROUP BY project_owner, project_project_name, DATE(closed_at) ORDER BY DATE(closed_at) ASC;

SELECT CONCAT(p.project_owner, '_', p.project_name) AS project, DATE(c.commit_date) AS time, COUNT(DISTINCT c.author) AS num_developers FROM commit c JOIN agraph ag ON c.a_graph_id = ag.id JOIN project p ON ag.project_project_owner = p.project_owner AND ag.project_project_name = p.project_name GROUP BY p.project_owner, p.project_name, DATE(c.commit_date) ORDER BY DATE(c.commit_date) ASC;

SELECT ag.project_project_owner AS repo_owner, ag.project_project_name AS repo_project_name, COUNT(DISTINCT c.author) AS num_developers FROM commit c JOIN agraph ag ON c.a_graph_id = ag.id GROUP BY ag.project_project_owner, ag.project_project_name ORDER BY num_developers DESC;

SELECT CONCAT(p.project_owner, '_', p.project_name) AS project, DATE(c.commit_date) AS time, SUM(fc.total_deletions) AS total_deletions FROM file_change fc JOIN commit_file_changes cfc ON fc.id = cfc.file_changes_id JOIN commit c ON cfc.commit_sha = c.sha JOIN agraph ag ON c.a_graph_id = ag.id JOIN project p ON ag.project_project_owner = p.project_owner AND ag.project_project_name = p.project_name GROUP BY p.project_owner, p.project_name, c.commit_date ORDER BY c.commit_date ASC;

SELECT CONCAT(p.project_owner, '_', p.project_name) AS project, DATE(c.commit_date) AS time, SUM(fc.total_changes) AS total_changes FROM file_change fc JOIN commit_file_changes cfc ON fc.id = cfc.file_changes_id JOIN commit c ON cfc.commit_sha = c.sha JOIN agraph ag ON c.a_graph_id = ag.id JOIN project p ON ag.project_project_owner = p.project_owner AND ag.project_project_name = p.project_name GROUP BY p.project_owner, p.project_name, c.commit_date ORDER BY c.commit_date ASC;

SELECT CONCAT(p.project_owner, '_', p.project_name) AS project, DATE(c.commit_date) AS time, SUM(fc.total_additions) AS total_additions FROM file_change fc JOIN commit_file_changes cfc ON fc.id = cfc.file_changes_id JOIN commit c ON cfc.commit_sha = c.sha JOIN agraph ag ON c.a_graph_id = ag.id JOIN project p ON ag.project_project_owner = p.project_owner AND ag.project_project_name = p.project_name GROUP BY p.project_owner, p.project_name, c.commit_date ORDER BY c.commit_date ASC;

SELECT CONCAT(p.project_owner, '_', p.project_name) AS project, DATE(c.commit_date) AS time, COUNT(c.sha) AS commitCount FROM commit c JOIN agraph ag ON c.a_graph_id = ag.id JOIN project p ON p.project_owner = ag.project_project_owner AND p.project_name = ag.project_project_name GROUP BY p.project_owner, p.project_name, time ORDER BY time ASC;

Pull Request Review Time:
SELECT CONCAT(project_owner, '_', project_project_name) AS project, DATE(created_at) AS time, AVG(TIMESTAMPDIFF(DAY, created_at, closed_at)) AS avg_review_time FROM project_pull where project_pull.closed_at IS NOT NULL GROUP BY project_owner, project_project_name, time ORDER BY time ASC;
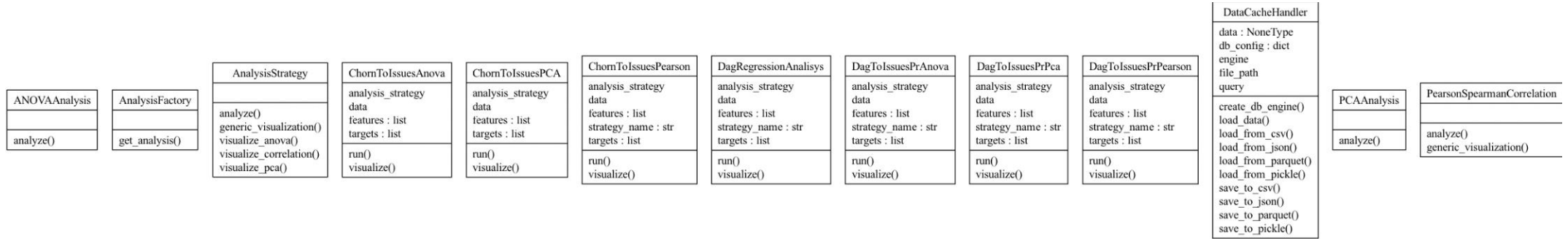
SELECT CONCAT(project_owner, ':', project_project_name) AS project, DATE(created_at) AS time, AVG(TIMESTAMPDIFF(DAY, created_at, closed_at)) AS avg_review_time FROM project_issue where project_issue.closed_at IS NOT NULL GROUP BY project_owner, project_project_name, time ORDER BY time ASC;
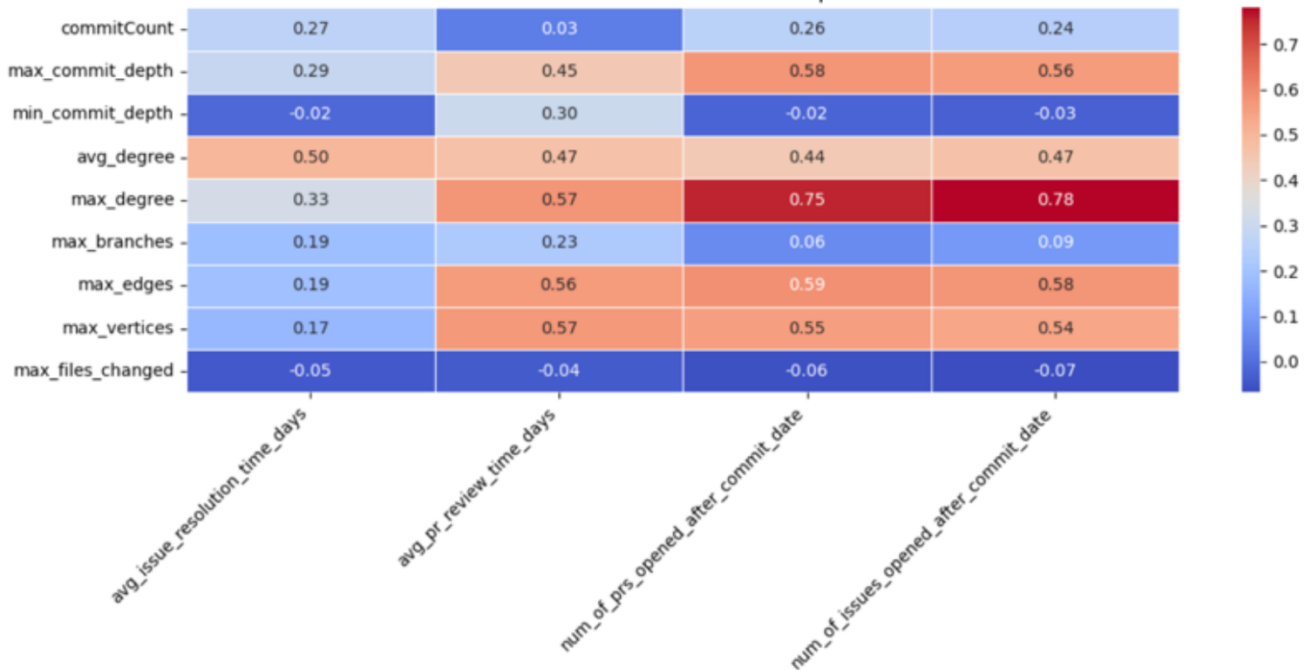
# Reporting

# Data Analisys

**ANOVAAnalysis**

analyze()

---

**AnalysisFactory**

get_analysis()

---

**AnalysisStrategy**

analyze()
generic_visualization()
visualize_anova()
visualize_correlation()
visualize_pca()

---

**ChornToIssuesAnova**

analysis_strategy
data
features : list
targets : list

run()
visualize()

---

**ChornToIssuesPCA**

analysis_strategy
data
features : list
targets : list

run()
visualize()

---

**ChornToIssuesPearson**

analysis_strategy
data
features : list
strategy_name : str
targets : list

run()
visualize()

---

**DagRegressionAnalisys**

analysis_strategy
data
features : list
strategy_name : str
targets : list

run()
visualize()

---

**DagToIssuesPrAnova**

analysis_strategy
data
features : list
strategy_name : str
targets : list

run()
visualize()

---

**DagToIssuesPrPca**

analysis_strategy
data
features : list
strategy_name : str
targets : list

run()
visualize()

---

**DagToIssuesPrPearson**

analysis_strategy
data
features : list
strategy_name : str
targets : list

run()
visualize()

---

**DataCacheHandler**

data : NoneType
db_config : dict
engine
file_path
query

create_db_engine()
load_data()
load_from_csv()
load_from_json()
load_from_parquet()
load_from_pickle()
save_to_csv()
save_to_json()
save_to_parquet()
save_to_pickle()

---

**PCAAnalysis**

analyze()

---

**PearsonSpearmanCorrelation**

analyze()
generic_visualization()

Pearson Correlation Heatmap


Spearman Correlation Heatmap

Higher complexity in commit structures (max_commit_depth, max_degree, avg_degree) is strongly associated with longer issue resolution times and higher PR review times.

More interconnected commit graphs (higher max_degree) correspond with more PRs and issues being opened after commit dates.

The number of branches (max_branches) and number of edges (max_edges) are also positively correlated with issue and PR-related metrics.

Next steps:

Research relations:
- Does the degree of a node (number of merges) correlate with the number of errors in actions in a CI pipeline?
- Does the number of branches commit went through correlate with the number of issues it could cause?
- Do commits with more than one parent (merge commits) have a significantly different number of PRs compared to linear commits?
- Do projects with different numbers of branches have different CI pipeline execution errors?
- What are the main factors influencing repository complexity? (e.g., branches, PRs, actions)
- Try to group repositories based on their commit structure