

IMT 542



GitHub User Analysis Project

Name: Xinyi (Eve) Wang

Overview

Info story or who is the user and why is it important to them

Existing structure and FAIR assessment

How you decided to improve the structure

What your new structure is

How would quality be identified and addressed



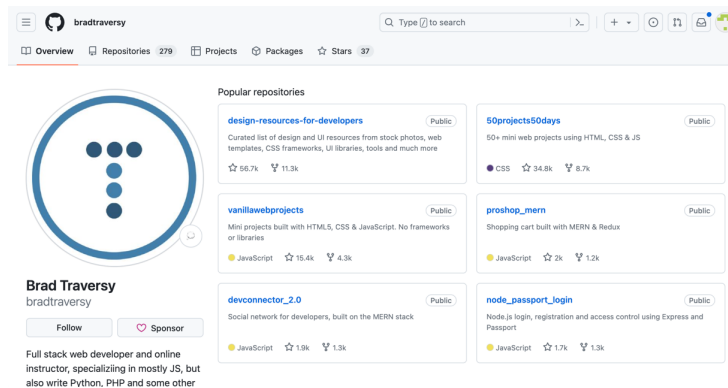
Info story

Recruiters - > Candidate's GitHub

- > in 2 minutes, determine: Are you a proficient Java developer? How good are you?

Problem:

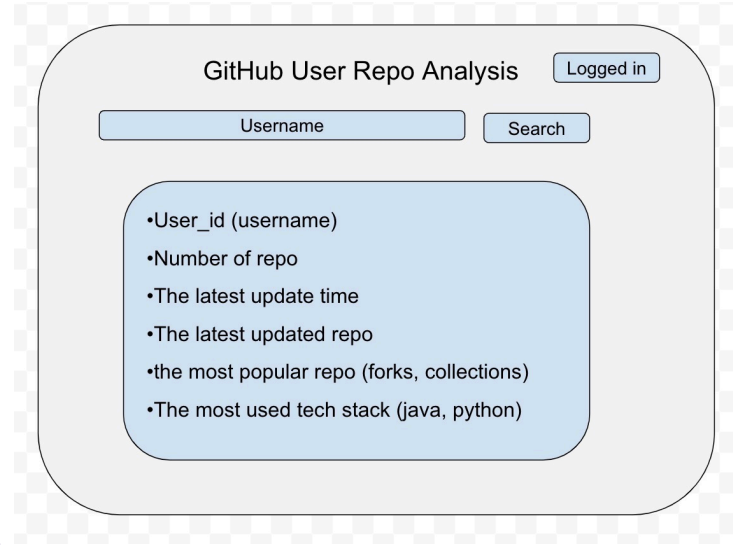
- Too much information
- What is their most frequently used language?
- Difficult to generate insights on candidates
- in a quick and uniform manner



Info story

Solution: Automated GitHub user analysis

- Retrieving data available on the GitHub webpage
- Filtering useful information for recruiters
- Analyze and add new fields with additional insights
- Summarize the user and repo insights
- Help recruiters find candidates with the right skillsets



Information wireframe



Existing structure and FAIR assessment

Existing structure from GitHub web page

The screenshot shows the GitHub profile of user EveWang. The profile includes a green and white pixelated avatar, the name 'EveWang', and the username 'EveWangUW'. It shows 0 followers and 2 following. The 'Popular repositories' section lists four public repositories: BookStoreApp (JavaScript, 1 star), MovieReviewApp (JavaScript), CDK-S3-Project (TypeScript, described as 'cdk project with 2 s3 buckets'), HackerNewsIOSApp (Swift), Hangman-Game (TypeScript), and CICD-Pipeline-with-CDK (TypeScript). A '145 contributions in the last year' heatmap is displayed, showing activity from June 2023 to May 2024. The 'Contribution settings' are set to '2024'.

The screenshot shows the repository page for 'IMT542_GitHubUserAnalysis_Xinyi-Eve-Wang'. The repository is public and has 12 commits. The file list includes: build, node_modules, public, src, .DS_Store, IM_README.md, IMT 542 Final pre.pptx, LICENSE, README.md, package-lock.json, package.json, readme_08_Xinyi Wang.md, and tsconfig.json. The 'About' section states: 'No description, website, or topics provided.' The 'Releases' section states: 'No releases published. Create a new release.' The 'Packages' section states: 'No packages published. Publish your first package.' The 'Languages' section shows a bar chart with TypeScript at 94.7% and HTML at 5.3%. The 'Suggested workflows' section is empty.



Existing structure and FAIR assessment

Existing structure from
GitHub REST API endpoints

1. User information

```
"login": "EveWangUW",
"id": 142384748,
"node_id": "U_kgDOCHyebA",
"avatar_url": "https://avatars.githubusercontent.com/u/142384748?v=4",
"gravatar_id": "",
"url": "https://api.github.com/users/EveWangUW",
"html_url": "https://github.com/EveWangUW",
"followers_url": "https://api.github.com/users/EveWangUW/followers",
"following_url": "https://api.github.com/users/EveWangUW/following{/other_user}",
"gists_url": "https://api.github.com/users/EveWangUW/gists{/gist_id}",
"starred_url": "https://api.github.com/users/EveWangUW/starred{/owner}/{repo}",
"subscriptions_url": "https://api.github.com/users/EveWangUW/subscriptions",
"organizations_url": "https://api.github.com/users/EveWangUW/orgs",
"repos_url": "https://api.github.com/users/EveWangUW/repos",
"events_url": "https://api.github.com/users/EveWangUW/events{/privacy}",
"received_events_url": "https://api.github.com/users/EveWangUW/received_events",
"type": "User",
"site_admin": false,
"name": "EveWang",
"company": null,
"blog": "",
"location": null,
"email": null,
"hireable": null,
"bio": null,
"twitter_username": null,
"public_repos": 25,
"public_gists": 0,
"followers": 0,
"following": 2,
"created_at": "2023-08-16T07:12:54Z",
"updated_at": "2024-05-03T01:02:37Z"
```



Existing structure and FAIR assessment

Existing structure from
GitHub REST API endpoints
2. User repo information

```
{
  "id": 793294266,
  "node_id": "R_kgDOL0ilug",
  "name": "ActiveRecall-StudyBestFriend",
  "full_name": "EveWangUM/ActiveRecall-StudyBestFriend",
  "private": false,
  "owner": {
    "login": "EveWangUM",
    "id": 142384748,
    "node_id": "U_kgD0ChyebA",
    "avatar_url": "https://avatars.githubusercontent.com/u/142384748?v=4",
    "gravatar_id": "",
    "url": "https://api.github.com/users/EveWangUM",
    "html_url": "https://github.com/EveWangUM",
    "followers_url": "https://api.github.com/users/EveWangUM/followers",
    "following_url": "https://api.github.com/users/EveWangUM/following{/other_user}",
    "gists_url": "https://api.github.com/users/EveWangUM/gists{/gist_id}",
    "starred_url": "https://api.github.com/users/EveWangUM/starred{/owner}/{repo}",
    "subscriptions_url": "https://api.github.com/users/EveWangUM/subscriptions",
    "organizations_url": "https://api.github.com/users/EveWangUM/orgs",
    "repos_url": "https://api.github.com/users/EveWangUM/repos",
    "events_url": "https://api.github.com/users/EveWangUM/events{/privacy}",
    "received_events_url": "https://api.github.com/users/EveWangUM/received_events",
    "type": "User",
    "site_admin": false
  },
  "html_url": "https://github.com/EveWangUM/ActiveRecall-StudyBestFriend",
  "description": null,
  "fork": false,
  "url": "https://api.github.com/repos/EveWangUM/ActiveRecall-StudyBestFriend",
  "forks_url": "https://api.github.com/repos/EveWangUM/ActiveRecall-StudyBestFriend/forks",
  "keys_url": "https://api.github.com/repos/EveWangUM/ActiveRecall-StudyBestFriend/keys{/key_id}",
  "collaborators_url": "https://api.github.com/repos/EveWangUM/ActiveRecall-StudyBestFriend/collaborators{/collaborator}",
  "teams_url": "https://api.github.com/repos/EveWangUM/ActiveRecall-StudyBestFriend/teams",
  "hooks_url": "https://api.github.com/repos/EveWangUM/ActiveRecall-StudyBestFriend/hooks",
  "issue_events_url": "https://api.github.com/repos/EveWangUM/ActiveRecall-StudyBestFriend/issues/events{/number}",
  "events_url": "https://api.github.com/repos/EveWangUM/ActiveRecall-StudyBestFriend/events",
  "assignees_url": "https://api.github.com/repos/EveWangUM/ActiveRecall-StudyBestFriend/assignees{/user}",
  "branches_url": "https://api.github.com/repos/EveWangUM/ActiveRecall-StudyBestFriend/branches{/branch}",
  "tags_url": "https://api.github.com/repos/EveWangUM/ActiveRecall-StudyBestFriend/tags",
  "blobs_url": "https://api.github.com/repos/EveWangUM/ActiveRecall-StudyBestFriend/git/blobs{/sha}",
  "git_tags_url": "https://api.github.com/repos/EveWangUM/ActiveRecall-StudyBestFriend/git/tags{/sha}",
  "git_refs_url": "https://api.github.com/repos/EveWangUM/ActiveRecall-StudyBestFriend/git/refs{/sha}",
  "trees_url": "https://api.github.com/repos/EveWangUM/ActiveRecall-StudyBestFriend/git/trees{/sha}",
  "statuses_url": "https://api.github.com/repos/EveWangUM/ActiveRecall-StudyBestFriend/statuses{/sha}",
  "languages_url": "https://api.github.com/repos/EveWangUM/ActiveRecall-StudyBestFriend/languages",
  "stargazers_url": "https://api.github.com/repos/EveWangUM/ActiveRecall-StudyBestFriend/stargazers",
  "contributors_url": "https://api.github.com/repos/EveWangUM/ActiveRecall-StudyBestFriend/contributors",
  "subscribers_url": "https://api.github.com/repos/EveWangUM/ActiveRecall-StudyBestFriend/subscribers",
  "subscription_url": "https://api.github.com/repos/EveWangUM/ActiveRecall-StudyBestFriend/subscription",
  "commits_url": "https://api.github.com/repos/EveWangUM/ActiveRecall-StudyBestFriend/commits{/sha}",
  "git_commits_url": "https://api.github.com/repos/EveWangUM/ActiveRecall-StudyBestFriend/git/commits{/sha}",
  "comments_url": "https://api.github.com/repos/EveWangUM/ActiveRecall-StudyBestFriend/comments{/number}",
  "issue_comment_url": "https://api.github.com/repos/EveWangUM/ActiveRecall-StudyBestFriend/issues/comments{/number}"
}
```



Existing structure and FAIR assessment

Existing structure from
GitHub REST API endpoints

2. User repo information

```
"issue_comment_url": "https://api.github.com/repos/EveWangUM/ActiveRecall-StudyBestFriend/issues/comments/{number}",
"contents_url": "https://api.github.com/repos/EveWangUM/ActiveRecall-StudyBestFriend/contents/{+path}",
"compare_url": "https://api.github.com/repos/EveWangUM/ActiveRecall-StudyBestFriend/compare/{base}...{head}",
"merges_url": "https://api.github.com/repos/EveWangUM/ActiveRecall-StudyBestFriend/merges",
"archive_url": "https://api.github.com/repos/EveWangUM/ActiveRecall-StudyBestFriend/{archive_format}/{ref}",
"downloads_url": "https://api.github.com/repos/EveWangUM/ActiveRecall-StudyBestFriend/downloads",
"issues_url": "https://api.github.com/repos/EveWangUM/ActiveRecall-StudyBestFriend/issues/{number}",
"pulls_url": "https://api.github.com/repos/EveWangUM/ActiveRecall-StudyBestFriend/pulls/{number}",
"milestones_url": "https://api.github.com/repos/EveWangUM/ActiveRecall-StudyBestFriend/milestones/{number}",
"notifications_url": "https://api.github.com/repos/EveWangUM/ActiveRecall-StudyBestFriend/notifications?since=all,participating",
"labels_url": "https://api.github.com/repos/EveWangUM/ActiveRecall-StudyBestFriend/labels/{name}",
"releases_url": "https://api.github.com/repos/EveWangUM/ActiveRecall-StudyBestFriend/releases/{id}",
"deployments_url": "https://api.github.com/repos/EveWangUM/ActiveRecall-StudyBestFriend/deployments",
"created_at": "2024-04-28T23:51:14Z",
"updated_at": "2024-04-28T23:52:01Z",
"pushed_at": "2024-04-28T23:51:24Z",
"git_url": "git://github.com/EveWangUM/ActiveRecall-StudyBestFriend.git",
"ssh_url": "git@github.com:EveWangUM/ActiveRecall-StudyBestFriend.git",
"clone_url": "https://github.com/EveWangUM/ActiveRecall-StudyBestFriend.git",
"svn_url": "https://github.com/EveWangUM/ActiveRecall-StudyBestFriend",
"homepage": null,
"size": 358,
"stargazers_count": 0,
"watchers_count": 0,
"language": "JavaScript",
"has_issues": true,
"has_projects": true,
"has_downloads": true,
"has_wiki": true,
"has_pages": false,
"has_discussions": false,
"forks_count": 0,
"mirror_url": null,
"archived": false,
"disabled": false,
"open_issues_count": 0,
"license": null,
"allow_forking": true,
"is_template": false,
"web_commit_signoff_required": false,
"topics": [],
"visibility": "public",
"forks": 0,
"open_issues": 0,
"watchers": 0,
"default_branch": "main",
},
{id": 777993033,
"node_id": "R_kgDOL1875Q",
```


Existing structure and FAIR assessment

FAIR assessment - GOOD:

F1: Both the user data and the user repo data have globally unique identifiers such as login and id.

F3: The metadata clearly and explicitly include the identifier of the data they describe, for example, the repo data includes the repo owner's information

F2: There is also rich metadata on the user and repository details, such as 'updated_at', 'created_at', profile URLs, and owner details

A1, A1.1: The data are retrievable via standardized, open, and free APIs provided by GitHub.

I1: These data use the JSON language format, which is a formal, accessible, shared, and broadly applicable language for knowledge representation

I3: the various URLs included in the user and the repo data direct the users to related resources, for example, the fields 'followers' and 'commits' in the user and repo data include URLs to various resources such as 'followers_url' and 'commits_url'.



Existing structure and FAIR assessment

FAIR assessment – deficiencies that need to be improved

A2: if the user account or the repo is deleted, the metadata might not be accessible anymore, which may not conform to the A2 principle.

R1: In terms of reusability, although rich attributes are provided, some of the attributes are not relevant to our recruiters, such as 'avater_url', 'is_template', not conforming to R1

R1.1: In the user and user repo, most of the data lack license information, not conforming to the R1.1.

R1.2: Finally, although time related attributes such as 'updated_at', 'created_at' provide some provenance information, but more details could be included to achieve the R1.2 principle.



Existing structure and portability assessment

Portability assessment – deficiencies that needs to be improved

Clean meaning: raw JSON with too much info and lack of summary

Accessibility: look at the page directly or make API calls

Transparency: difficulty understanding the field

Interoperability: JSON format only, may not work well on mobile devices.

Usability: difficult to use directly for candidate assessment

Structure: lack of summary and structure



How you decided to improve the structure

To improve accessibility (A),

1. implement a backup system to retain metadata even if the user account or repo is deleted, which can be accomplished by allowing the user to **download the data** from the web application

To improve reusability (R),

1. **Remove unnecessary attributes** like 'avator_url' and focus on displaying information related to recruiters.

2. Add relevant **license** for this application

3. **Add more analytics details** and provenance information to provide recruiters with useful insights into candidate profiles.



How you decided to improve the structure

Portability improvement:

Clean meaning:

raw JSON with too much info and lack of summary -> summary of insights

Accessibility:

look at the page directly or make API calls -> interactive frontend to present the analysis result + provide downloadable data

Transparency:

difficulty understanding the field -> detailed explanation documentation



How you decided to improve the structure

Portability improvement:

Interoperability:

JSON format only, may not work well on mobile devices

-> typescript for typing and data schema,

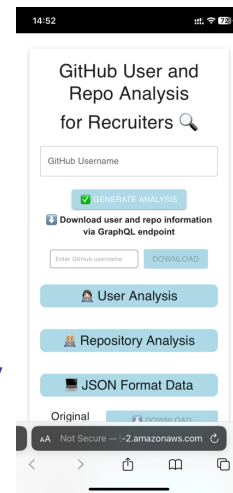
-> React for interoperability on web and mobile, more data format such as CSV

Usability:

difficult to use directly for candidate assessment -> analyze based on raw data and generate deeper insights on the user and their repos

Structure:

lack of summary and structure -> combine user and repo insights in one JSON file



What your new structure is

- > **GitHub User and Repo Analysis for Recruiters (hosted on Amazon S3)**
- > <http://githubanalysisave.s3-website-ap-southeast-2.amazonaws.com/>

GitHub User and Repo Analysis for Recruiters 🔍

✓ GENERATE ANALYSIS

Download user and repo information via GraphQL endpoint

Enter GitHub username

DOWNLOAD

User Analysis

Repository Analysis

JSON Format Data

Original User Data:

DOWNLOAD ORIGINAL USER DATA

Original Repo Data:

DOWNLOAD ORIGINAL REPOSITORY DATA

New User Data:

DOWNLOAD NEW USER DATA

New Repo Data:

DOWNLOAD NEW REPO DATA

New User and Repo Data:

DOWNLOAD COMBINED

🔧 Filtered User Repository Details (first 5 repos)



What your new structure is

- > Data manipulation and transformation
- > See code!



What your new structure is

Information itself:

Existing information structure:

too much information

New information structure:

1. Information has been trimmed

-> retain information that is useful to recruiters in analyzing a candidate's programming skill level.

2. Added new fields: e.g.: the time of the user on GitHub

-> analyzed and summarized information from the user account and all the repositories' information.



What your new structure is

Structure/format:

Existing information structure:

raw JSON data directly retrieved from the GitHub API, includes overall details fields that may not be relevant to the recruiters (e.g. node IDs or avatar URLs)

New information structure:

filtered and summarized data that includes only relevant fields (e.g. user name, repository name)

and additional computed fields (e.g. total time the user has been on GitHub and fields about their programming skill levels)

-> providing more relevant information to the recruiters.



What your new structure is

Access methodology:

Existing information structure:

- Directly access the GitHub webpage
- Manual API calls

New information structure:

- Web application hosted on Amazon S3
- GraphQL API calls
- REST API calls



How would quality be identified and addressed

Lasting change:

- Provide **downloadable** original and new information structure

Durable and Robust Application:

- Apply rigorous **unit-testing**, including:
 - Sending requests - > Hitting the GitHub API endpoint
 - Getting responses -> Rendering components on the frontend based on responses
- Use TypeScript unit-testing frameworks with a goal of achieving 99%+ test coverage.

Application and data security

- Use an **MIT license** for the code to ensure open-source accessibility.
- The application is hosted on **Amazon S3 bucket** for security



Future improvements

- > Storage of past user and repo performance (a detailed timeline of the GitHub user's performance)
- > More personalized GraphQL endpoint usage, allowing users to extract the data they want
- > Data security



Q & A

Thank you! :)