# SiteScanning Developer Handoff

## Key Repositories

The README files in the repositories listed below are comprehensive and provide all the necessary information for developers to begin working on these projects.

| Repository Name | Description | Link to Repo |
| --- | --- | --- |
| site-scanning | Serves as the home to the Site Scanning project. All project issues are managed through this repository. | https://github.com/GSA/site-scanning |
| site-scanning-documentation | This repository houses all of the documentation related to the Site Scanning project. This is the first place you should look if you are trying to find anything project related. | https://github.com/GSA/site-scanning-documentation |
| site-scanning-analysis | This repository houses numerous reports. Some of these reports are scheduled daily through actions and some are one-off datasets used by the federal-website-index. | https://github.com/GSA/site-scanning-analysis |
| federal-website-index | This is the home of the Federal Website Indexer application. The indexer runs daily and builds a list of all public federal government websites. | https://github.com/GSA/federal-website-index |
| site-scanning-snapshots | This repository stores a variety of one-off scan data. | https://github.com/GSA/site-scanning-snapshots |
| site-scanning-engine | This repository houses both the Site Scanning Engine and the API that is used to retrieve the resulting data from the Site Scanning Engine. | https://github.com/GSA/site-scanning-engine |

## Notes

Among the repositories listed above, developers should pay particular attention to the **site-scanning-engine**, **federal-website-index**, and **site-scanning-analysis** repositories. Each of these includes documentation on how to set up a local development environment to run the projects locally.

# Scheduled Workflows

This section outlines the GitHub Actions implemented to ensure the project runs smoothly and on autopilot. The actions are categorized by repository, with descriptions, schedules, and links for each.

## site-scanning-engine

| Action Name | Description | Schedule | Link to Action |
|---|---|---|---|
| Restart scan worker/consumer | Restarts the scan consumer in cloud.gov to prepare for scan | Runs daily at 00:00 UTC | Link to Action: site-scanning-engine/restart -worker.yml |
| Enqueue Scans | Queues up the scans loaded from the ingest workflow | Runs at 00:00 UTC on Mon / Wed / Fri | Link to Action: site-scanning-engine/enque ue-scans.yml |
| Create S3 snapshot | Generates the legacy version of the s3 snapshots | Runs daily at 12:00 UTC | Link to Action: site-scanning-engine/create -snapshot.yml |
| Requeue stale scans | Queues up the urls that have not had scan data updated for over 3 days. | Runs at 12:00 UTC on Tue / Thur / Sat | Link to Action: site-scanning-engine/reque ue-stale-scans.yml |
| Create daily S3 snapshots | Generates the snapshots located in S3 and archives older versions | Runs daily at 12:15 UTC | Link to Action: site-scanning-engine/create -daily-snapshots.yml |

| Action Name | Description | Schedule | Link to Action |
|---|---|---|---|
| Create a11y S3 snapshot | Generates the a11y snapshots located in S3 | Runs daily at 12:15 UTC | Link to Action: site-scanning-engine/create-a11y-snapshot.yml |
| Ingest | Ingests all URLs from the federal-site-index | Runs daily at 22:15 UTC | Link to Action: site-scanning-engine/ingest.yml |

## federal-website-index

| Action Name | Description | Schedule | Link to Action |
|---|---|---|---|
| Build DAP Top 100,000 List | Builds the DAP Top 100,00 list that is consumed by the indexer | Runs daily at 21:00 UTC | Link to Action: federal-website-index/build-dap-top-list.yml |
| Build Final URL List | Builds the final-url list that is consumed by the indexer | Runs daily at 20:30 UTC | Link to Action: federal-website-index/build-finalurl-list.yml |
| Build target url list using TypeScript | Runs the primary federal website indexer that generates the url list that is consumed by the site-scanning-engine | Runs daily at 20:15 UTC | Link to Action: federal-website-index/build-list-js.yml |

## site-scanning-analysis

| Action Name | Description | Schedule | Link to Action |
|---|---|---|---|
| Generate all reports | Runs all reports within the analysis repo that are not individually scheduled | Runs daily at 13:00 UTC | Link to Action: site-scanning-analysis/generate-all-reports.yml |
| Generate IDEA Reports | Runs the IDEA Reports | Runs at 00:00 UTC on Thursdays | Link to Action: site-scanning-analysis/generate-idea-reports.yml |
| Generate Unique Website List | Runs the unique-website-list report | Runs daily at 14:00 UTC | Link to Action: site-scanning-analysis/generate-unique-website-list.yml |
| Generate URLs Missing From Snapshot | Generates a report of URLs that exist in the index file, but do not have an entry in the daily snapshot | Runs daily at 21:00 UTC | Link to Action: site-scanning-analysis/generate-urls-missing-from-snapshot.yml |
| Run Smoke Tests | Runs the smoke tests | Runs daily at 13:30 UTC | Link to Action: site-scanning-analysis/smoke-tests.yml |

## Site Scanning Engine/API: Development and Deployment Overview

The **README** file in the Site Scanning Engine repository provides clear and thorough instructions for running both the scan engine and the API in a local development environment.

### Development Workflow

When you're ready to deploy changes, follow these steps:

- **Create a Pull Request (PR):**
  Open a PR to merge your feature branch into the main branch.

- **Automated Checks:**
  Upon creating the PR, several GitHub Actions workflows will run:

  - **Unit and integration tests** to ensure code stability.

  - **Test deployment**, which deploys the scan engine and API to the **development environment** to verify that the deployment process completes successfully.

- **Merge to Main:**
  Once all checks pass, you can safely merge your branch into **main**.

- **Production Deployment:**
  After merging, manually trigger the **"Deploy" GitHub Action** to deploy the scan engine and API to the **production environment**.

### Logging and Monitoring

The **Scan Engine** generates extensive logging data, which can be accessed via logs.fr.cloud.gov. All logs are stored in **OpenSearch**, providing powerful search and filtering capabilities.

A number of saved searches are available to help streamline log analysis. These saved queries are all prefixed with **sse** for easy identification.

# Useful Links

| Name | Description | Link |
|---|---|---|
| Access the API | Full instructions on how the we can access the API in the production environment | https://open.gsa.gov/api/site-scanning-api/ |
| Access the Data | Outlines the numerous ways the data can be accessed | https://digital.gov/guides/site-scanning/data/ |

| Name | Description | Link |
|------|-------------|------|
| Architecture Documentation | Contains a diagram that outlines the infrastructure of the scan-engine application | https://github.com/GSA/site-scanning-engine/blob/main/docs/architecture/README.md |
| Data Dictionary | Defines each field within the scan-engine snapshot | https://github.com/GSA/site-scanning-documentation/blob/main/data/Site_Scanning_Data_Dictionary.csv |
| How the Federal Website Index is created | This lays out the steps involved in building the Federal Website Index | https://github.com/GSA/federal-website-index/blob/main/process/index-creation.md |
| How the Scans Run | This lays out the different scans the site-scanning-engine runs along with where they are located | https://github.com/gsa/site-scanning-documentation?tab=readme-ov-file#understanding-the-data |
| Scan statuses explained | Explains the different scan statuses that returned from the scan-engine | https://github.com/GSA/site-scanning-documentation/blob/main/pages/scan_statuses.md |
| Site Scanning Issues Board | This is the main board for all issues for the Site Scanning project | https://github.com/GSA/site-scanning/issues |
| Technical Details | If in doubt, look here. | https://digital.gov/guides/site-scanning/technical-details |