# 4 - DataFlow Overview // Onboarding

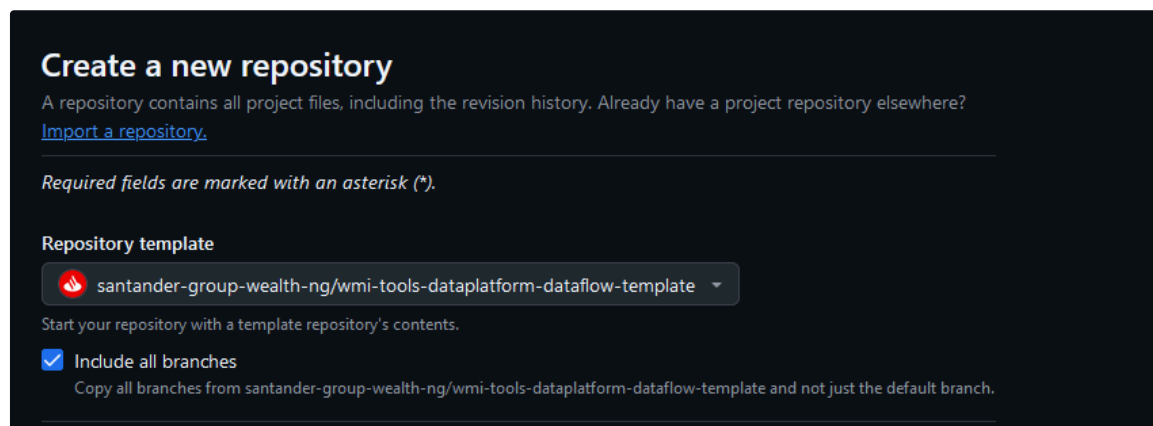## Introduction 🔗

## How do I generate the repository? 🔗

The repository should be generated using the following repository as a template  https://github.com/santander-group-wealth-ng/wmi-tools-dataplatform-dataflow-template  Connect your Github account
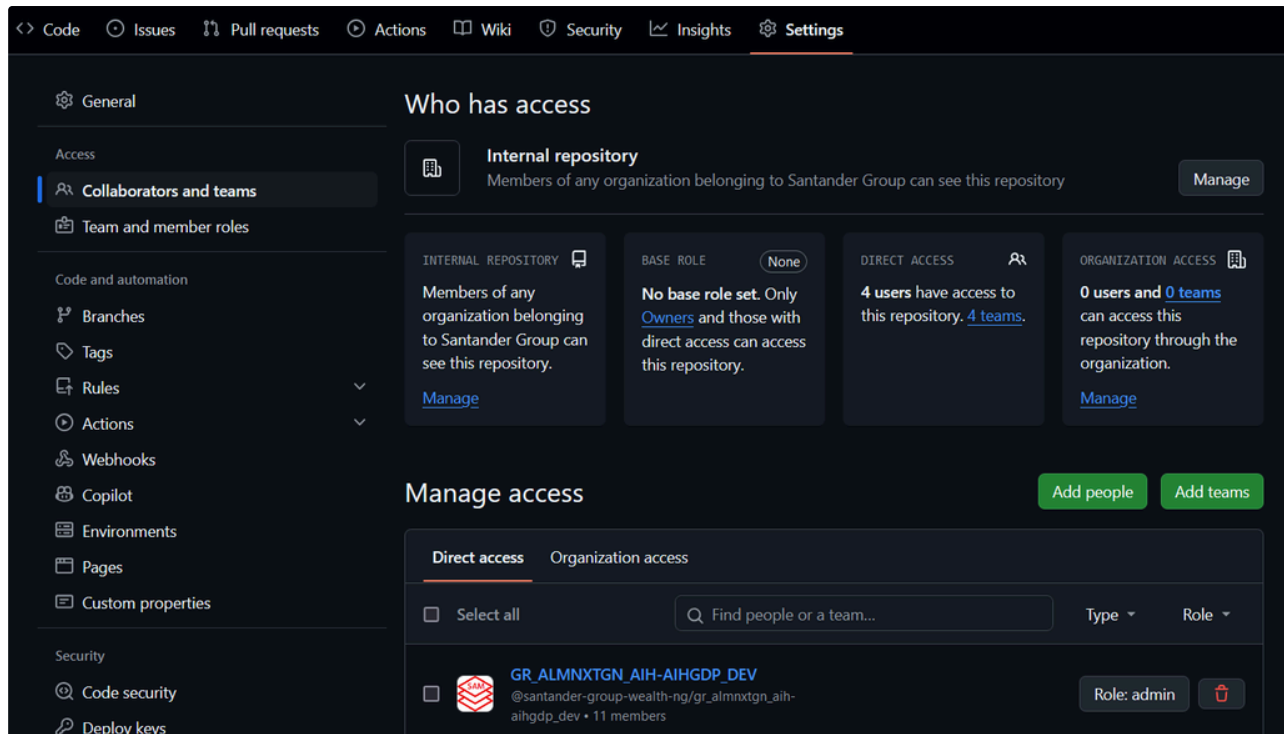
The steps to follow are as follows:

- Click on the `Use this template` button

- Important enable `include all branches`
- The name that should be given to the repository must follow the following format
  - aih-dataplatform-python-dataflow-{project name}_{data flow name} → aih-dataplatform-python-dataflow-morningstar_classes

⚠️ The user which done this task should remove his nominative permission and should add the GR_ALMNXTGN_AIH-AIHGDP_DEV team with admin role.

- The way for changing the permissions of the repository is going to "Settings" → "Collaborators and teams" → "Manage access". First, the technician should add the GR_ALMNXTGN_AIH-AIHGDP_DEV team by clicking on the button "Add teams" and give the admin role. Then, he should remove his self's nominative permission.



## Repository Structure 🔗

Once the **dataflow repository** is created and configured we will have a directory structure like this one:

```
Proyect/
├── .github/
├── data_catalog/
│   └── data_catalog.json
├── data_quality/
│   ├── rule1.gdq
│   └── rule2.gdq
├── test/
│   ├── files/
│   │   └── file1.txt
│   └── test_1.py
├── transformation/
```

```
|   ├── src/
|   |   ├── dataflow/
|   |   |─── default ├── raw_to_staging.py
|   |   |                 ├── staging_to_common.py
|   |   ├── utils/
|   |   └── tests/
|   ├── transformation.json
```

The **.github** folder in a newly created repository typically contains configuration files related to repository management and automation. Here's what it may include:

- `workflows/`:
  - Contains YAML files that define GitHub Actions workflows for CI/CD or other automation processes.
  - You can configure action folders to organize required and auxiliary actions more effectively.

> ⚠ This folder is synchronized with the template and cannot be changed its content.

## data_catalog folder 🔗

The **data_catalog** folder contains a `data_catalog.json` file that is a structured JSON document used to describe and manage the datasets within a project. It acts as a central repository of metadata, providing detailed information about the datasets, their structure, sources, and relationships. This file is particularly useful for organizing, documenting, and automating data workflows.

Dictionary Schema

`data_catalog.json`

```
 1   note: (1) Appears one time
 2         (n) Appears n times
 3   [
 4     (n) {
 5       (1)"name": "<table_name>",
 6       (1)"db": "<database_name>",
 7       (0,1)"external": "<value>",
 8       (1) "columns": [
 9         {
10           (n)"FieldName": "<type>"
11         }
12       ],
13       (1) "partitions": [
14         (n){
15             (1) "PartitionKey": "<partition_key>",
16             (0,1) "Type": "<data_type>"
17         }
18       ],
19       (0,1) "format": "<format_type>",
20       (0,1) "extra_sentences": "<additional_configurations>",
21       (1) "table_properties": "<value>",
22       (1) "LfTags": [
23         (n){
24           (1)"TagKey": "<tag>",
25           (1) "TagValues": [
26             (n)"<value>"
27
```

```
28          ]
29        }
30      ]
31    }
32  ]
33
34
```

| KEY 🔗 | DEFINITION 🔗 | OPTIONS 🔗 | NOTE 🔗 |
|---|---|---|---|
| name | Table Name | | The name should have the words separeted by underscore, for example, "sells_rips_prod". |
| external | Table Type | `EXTERNAL` | Not applicable to Iceberg tables, such as those in staging or common layers |
| db | Database Name | `<db>.<name>` | The name should have the words separeted by underscore, and the name of the layer at the end, for example, "funds_raw". |
| columns | Column Definitions | `(column_name column_type, ...)` | It should be a dictionary with the keys with the name of the column and the value with the type of the fields. |
| partitions | Partition Keys | `PARTITIONED BY (partition_column column_type)` | Iceberg tables do not include the type because it is defined in the column tags. Athena tables, on the other hand, include the type, but the field is not defined within the column tags. |
| format | Storage Format | `STORED AS <format>` (e.g., `TEXTFILE`, `PARQUET`) | Not applicable to Iceberg tables, such as those in staging or common layers |
| extra_sentences | Row Format | `ROW FORMAT DELIMITED FIELDS TERMINATED BY '<delimiter>'` | Not applicable to Iceberg tables, such as those in staging or common layers |
| table_properties | Table Properties | `TBLPROPERTIES ("key"="value")` | |
| LFTags | Tags (Lake Formation or Glue Catalog tags) | Not part of DDL (applied post-creation by step_function `addLFTagToResource`) | It's a list of dictionaries with name of the TagKey and a list of TagValues. |

Example:

⌄ Click here to expand

```
1  [
2    {
3      "name": "sells_rips_prod",
4      "db": "funds_raw",
5      "external": "EXTERNAL",
6      "columns": {
7        "SecId":"string",
```

```
 8          "PerformanceId":"string",
 9          "Date":"string",
10          "Deleted":"string",
11          "Unit_BAS":"string",
12          "Unit_USD":"string",
13          "Unit_EUR":"string",
14          "Unit_GBP":"string",
15          "Unit_CHF":"string"
16        },
17        "partitions":[
18          {
19            "PartitionKey": "DataDate",
20            "Type": "string"
21          }
22        ],
23        "format": "TEXTFILE",
24        "extra_sentences": "ROW FORMAT DELIMITED FIELDS TERMINATED BY ';'",
25        "table_properties": "\"skip.header.line.count\"=\"1\"",
26        "LfTags": [
27          {
28            "TagKey": "domain",
29            "TagValues": [
30              "funds"
31            ]
32          },
33          {
34            "TagKey": "taxonomy",
35            "TagValues": [
36              "public"
37            ]
38          }
39        ]
40      },
41      {
42        "name": "sells_rips_correction",
43        "external": "EXTERNAL",
44        "db": "funds_raw",
45        "columns": {
46          "SecId":"string",
47          "PerformanceId":"string",
48          "Date":"string",
49          "Deleted":"string",
50          "Unit_BAS":"string",
51          "Unit_USD":"string",
52          "Unit_EUR":"string",
53          "Unit_GBP":"string",
54          "Unit_CHF":"string",
55          "Unit_DKK":"string",
56          "Unit_NOK":"string"
57        },
58        "partitions":[
59          {
60            "PartitionKey": "DataDate",
61            "Type": "string"
62          }
63        ],
64        "format": "TEXTFILE",
65        "extra_sentences": "ROW FORMAT DELIMITED FIELDS TERMINATED BY ';'",
```

```json
      "table_properties": "\"skip.header.line.count\"=\"1\"",
      "LfTags": [
        {
          "TagKey": "domain",
          "TagValues": [
            "funds"
          ]
        },
        {
          "TagKey": "taxonomy",
          "TagValues": [
            "public"
          ]
        }
      ]
    },
    {
      "name": "sells_rips",
      "db": "funds_staging",
      "columns": {
        "SecId":"string",
        "PerformanceId":"string",
        "Date":"date",
        "Deleted":"boolean",
        "Unit_BAS":"double",
        "Unit_USD":"double",
        "Unit_EUR":"double",
        "Unit_GBP":"double",
        "Unit_CHF":"double",
        "Unit_DKK":"double",
        "DataDate":"date"
      },
      "partitions":[
        {
          "PartitionKey": "DataDate"
        }
      ],
      "table_properties": "'table_type'='ICEBERG', 'format'='parquet', 'write_compression'='snappy'",
      "LfTags": [
        {
          "TagKey": "domain",
          "TagValues": [
            "funds"
          ]
        },
        {
          "TagKey": "taxonomy",
          "TagValues": [
            "public"
          ]
        }
      ]
    },
    {
      "name": "sells_rips_spark",
      "db": "funds_staging",
      "columns": {
        "SecId":"string",
```

```
124          "PerformanceId":"string",
125          "Date":"date",
126          "Deleted":"boolean",
127          "Unit_BAS":"double",
128          "Unit_USD":"double",
129          "Unit_EUR":"double",
130          "Unit_GBP":"double",
131          "Unit_CHF":"double",
132          "Unit_DKK":"double",
133          "Unit_NOK":"double",
134          "Unit_SEK":"double",
135          "Unit_JPY":"double",
136          "LastUpdate":"timestamp",
137          "Unit_SGD":"double",
138          "ReturnType":"int",
139          "Filled":"boolean",
140          "Unit_TWD":"double",
141          "Unit_HKD":"double",
142          "Unit_MYR":"double",
143          "Unit_CNY":"double",
144          "Unit_ILS":"double",
145          "Unit_INR":"double",
146          "Unit_CAD":"double",
147          "Unit_KWD":"double",
148          "Unit_PLN":"double",
149          "Unit_AUD":"double",
150          "Unit_THB":"double",
151          "Unit_KRW":"double",
152          "Unit_NZD":"double",
153          "DataDate": "date"
154        },
155        "partitions":[
156          {
157            "PartitionKey": "DataDate"
158          }
159        ],
160        "table_properties": "'table_type'='ICEBERG', 'format'='parquet', 'write_compression'='snappy'",
161        "LfTags": [
162          {
163            "TagKey": "domain",
164            "TagValues": [
165              "funds"
166            ]
167          },
168          {
169            "TagKey": "taxonomy",
170            "TagValues": [
171              "public"
172            ]
173          }
174        ]
175      },
176      {
177        "name": "sells_rips",
178        "db": "funds_common",
179        "columns": {
180          "SecId":"string",
181          "PerformanceId":"string",
```

```
182          "Date":"date",
183          "Deleted":"boolean",
184          "Unit_BAS":"double",
185          "Unit_USD":"double",
186          "Unit_EUR":"double",
187          "Unit_GBP":"double",
188          "Unit_CHF":"double",
189          "Unit_DKK":"double",
190          "Unit_NOK":"double",
191          "Unit_SEK":"double"
192        },
193        "partitions":[
194          {
195            "PartitionKey": "Date"
196          }
197        ],
198        "table_properties": "'table_type'='ICEBERG', 'format'='parquet', 'write_compression'='snappy'",
199        "LfTags": [
200          {
201            "TagKey": "domain",
202            "TagValues": [
203              "funds"
204            ]
205          },
206          {
207            "TagKey": "taxonomy",
208            "TagValues": [
209              "public"
210            ]
211          }
212        ]
213      },
214      {
215        "name": "sells_rips_spark",
216        "db": "funds_common",
217        "columns": {
218          "SecId":"string",
219          "PerformanceId":"string",
220          "Date":"date",
221          "Deleted":"boolean",
222          "Unit_BAS":"double",
223          "Unit_USD":"double",
224          "Unit_EUR":"double",
225          "DataDate": "date"
226        },
227        "partitions":[
228          {
229            "PartitionKey": "Date"
230          }
231        ],
232        "table_properties": "'table_type'='ICEBERG', 'format'='parquet', 'write_compression'='snappy'",
233        "LfTags": [
234          {
235            "TagKey": "domain",
236            "TagValues": [
237              "funds"
238            ]
239          },
```

```
240        {
241          "TagKey": "taxonomy",
242          "TagValues": [
243            "public"
244          ]
245        }
246      ]
247    }
248  ]
```

## data_quality folder 🔗

The **data_quality** folder contains different files that include the data quality rules written in Glue's Data Quality Language ( 📦 Data Quality Definition Language (DQDL) reference - AWS Glue ). Here is an example:

```
1  #Checks the primary key
2  IsPrimaryKey "performanceid" "date"
3  #Checks the column type
4  ColumnDataType "datadate" = "Date"
5  #Checks the format
6  CustomSql "SELECT COUNT(*) FROM primary WHERE CAST(dataDate AS STRING) REGEXP '^\d{4}-\d{2}-\d{2}$'" = 0
```

There is no separator between sentences. These sentences will be included in `data_catalog.json` each time there is a CI/CD action. Although it is not visible in the repository, it will be loaded into the AWS account.

For quality rules, a separate file will be created for each table and layer. The naming convention will be as follows: `funds_raw` for the raw layer, `funds_staging` for the staging layer, and `funds_common` for the common layer. After this, separated by a dot, the table name will follow, and then the file extension `.gdq` (Glue Data Quality).

Example of file name format:

{database name}.{table_name}.gdq → funds_common.morningstar_prices.gdq



The **test** folder contains unit test files along with their corresponding code. The developer can put the source data files used for the tests in the *files* subdirectory.

The **transformation** directory contains all the code required to transform data across the various layers of the data lake. A `transformation.json` file is a structured configuration used in data workflows used to define and document the rules and processes for transforming data between different layers in a data lake. It specifies transformation logic (including operations like filtering), renaming columns, and aggregating data. It also describes the output schema, data sources, destinations, and dependencies. The

**src/dataflows/default** directory contains the specific **python** code for transformations where the changes defined in `transformations.json` for each layer are executed. The utils directory should be used for helper functions or classes, and additional directories within `src/dataflows` for organizing other code as needed.

Structure of transformation.json

# Branching Model 🔗

We will adopt the **Gitflow** branching model, which follows these rules:



**Gitflow model**

**a. Main Branch (** `main` **):** 🔗

- Contains stable and production-ready code.
- Only updated with versions that have passed thorough testing.

**b. Development Branch (** `develop` **):** 🔗

- Contains the latest version of the code under development.
- Serves as the base for new features and team collaboration.

**c. Support Branches:** 🔗

**i. Feature Branches (** `feature` **):**

- Created from `develop`.
- Used to develop new features.
- Merged back into `develop` once complete.

**ii. Release Branches (** `release` **):**

- Created from `develop` when preparing a new version for production.
- Used for final adjustments, bug fixes, and testing.
- Merged into both `main` and `develop`.
- Example: `release/1.0`.

**iii. Hotfix Branches (** `hotfix` **):**

- Created from `main` to address critical issues in production.

- Merged into both `main` and `develop` .
- Sometimes, although it is rare, hotfixes can be applied to release branches.

# 3. CI/CD 🔗

**CI/CD (Continuous Integration/Continuous Deployment)** is a software development practice that automates the process of integrating code changes, testing, and deploying applications.

Which kind of action trigger ci/cd?

These are the most common reasons:

- **Code Changes:**
    - Pushing code to a repository.
    - Merging a pull request.
- **Manual Triggers:**
    - Workflow dispatch or manually starting a pipeline.

The most common trigger for this is a push action. It's worth emphasizing that the CI/CD process identifies the modified code and deploys only the affected parts. For example, if `data_catalog.json` is modified, only the CI/CD process related to the data catalog will be triggered.

Depending on your location, the trigger is executed in its corresponding environment. So if you are in `develop` , it will run the workflow in the `develop` repository; if you are in a feature, it will run in its respective repository; or if you are in `main` , it will run in the `main` repository.

It automatically upload and deploy catalog, transformations and orchestration processes.

If you want to check the workflow you will have to go to the **Actions** section in GitHub:

You can click on each job in the workflow to view the corresponding log:



It is also possible to trigger a workflow without performing any action on the repository, from **"deploy changes"**, with the option to choose what you want to launch and from which branch:

In order to "Run workflow" button appear, the workflow must to have been runned in the default branch of the repository. By default the main branch is the default branch, so, if a commit hasn't been done in the main branch (main), the button will not be visible. You should go to "Settings" → "General" → "Default branch" → "Switch to another branch" and select "develop" for make "Run workflow" button visible.
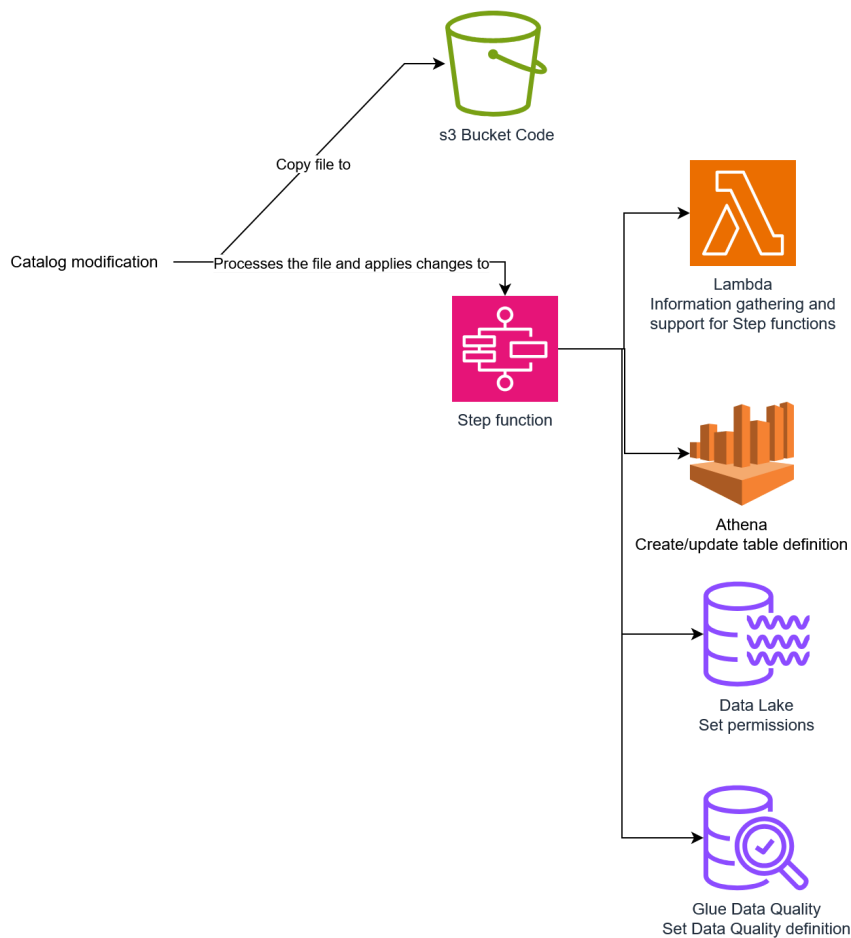


Default branch configuration menu

It is important to note that the CI/CD pipeline includes a step for **Fortify**, a tool used to analyze source code for vulnerabilities and security issues. This process scans the dataflow code using Fortify to detect potential security problems. Additionally, the pipeline includes a process for **Sonar**, a tool that assesses code quality and identifies issues such as bugs, vulnerabilities, and code smells. It is important to highlight that in the **dev** environment, passing Fortify and Sonar is not mandatory, but it is required in the **pre** and **pro** environments.
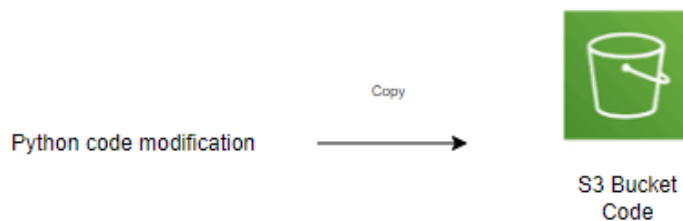
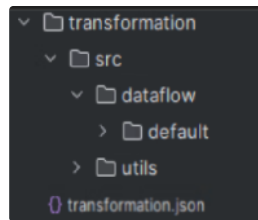There are three steps that the CI/CD process is prepared to deploy:

# Catalog 🔗

Each modification to the data_catalog.json file triggers two actions. The first one copies the code to its corresponding bucket, and the second calls the step function that applies the changes to the database and Glue Data Catalog (updating metadata, tables and quality rules)

## Python transformation code 🔗



Any change in the Python code within the transformation directory requires uploading all the transformation code, compressed, to the S3 bucket.

## Transformation.json 🔗



This process automates the update of **AWS Step Functions** by generating new `.asl.json` definitions based on a `transformations.json` file, storing them in an S3 bucket, and then using them to update the corresponding Step Functions. Each `.asl.json` file describes the workflow logic for a specific layer.