

Back Market Data Engineering Serve team Interview

Solution for the last question of the interview :

Eventually, more and more feature teams want to expose data and we want to provide a standard way of doing so. How would you build a tooling library and make it available to them (deployment, versioning, security, legal, etc.)?

A tooling library could be developed by the team responsible of the DWH to make external **tech** teams able to push data to the DWH without having to develop specific code and handle access control by themselves.

The tooling library could be designed like this :

- Available in a git repository as python code.
- Be pushed by the CI to a pypi index as a versioned package : every time the DWH team pushes changes on the code & update the version, a new tag would be generated and a new package available in pypi. This can be done easily by plugging the git repository to a CI/CD tool as Jenkins, Travis, or Gitlab CI.
- Handle access control through simple setup on their environment & with appropriate IAM roles given to their service account (created in Google Cloud IAM). A good way to do this is to define a separate service account by data producers group (eg FT) to only make them able to have access to data they are allowed to read, and write in datasets they can legitimately write.

The drawback of providing a simple “tooling library” for external teams is the team in charge of the DWH can’t easily have an eye on everything which is loaded, and the other team can’t take advantage of all the aspects around it which can be setup in the data platform (alerting, monitoring).

It can be mitigated by pushing some metrics directly in the library, for example by registering data in a data catalog or pushing metrics in a monitoring service directly through the library.

Other possible alternatives to make external teams able to loaded data could be to :

- Make it possible to create tables through external resources upload.
For example, an ETL (let’s say an Airflow DAG) which polls a GCS bucket to retrieve new data files uploaded with either a good nomenclature to create the appropriate tables :
`gcs://backmarket-external-teams/ftA/<table_name>_<date>.csv`
Or configuration files in a repository/in the GCS bucket to load the data in the right tables
- Use external tables with data hosted in GCS
- Make the FTs able to deploy ETL jobs on the data platform, either with a ready-to-use DAG (I’m still thinking with *Airflow*) factory, or top-level configuration files to define DAGs
- Provide a MS as intermediary to handle the loading (either by posting directly some data, or file deposit on a cloud file system + notification on the MS)