

Scalable Development of Event Driven Airflow DAGs

Subramanian Vellaiyan
Data Engineering



Ipsa Trivedi
Data Engineering



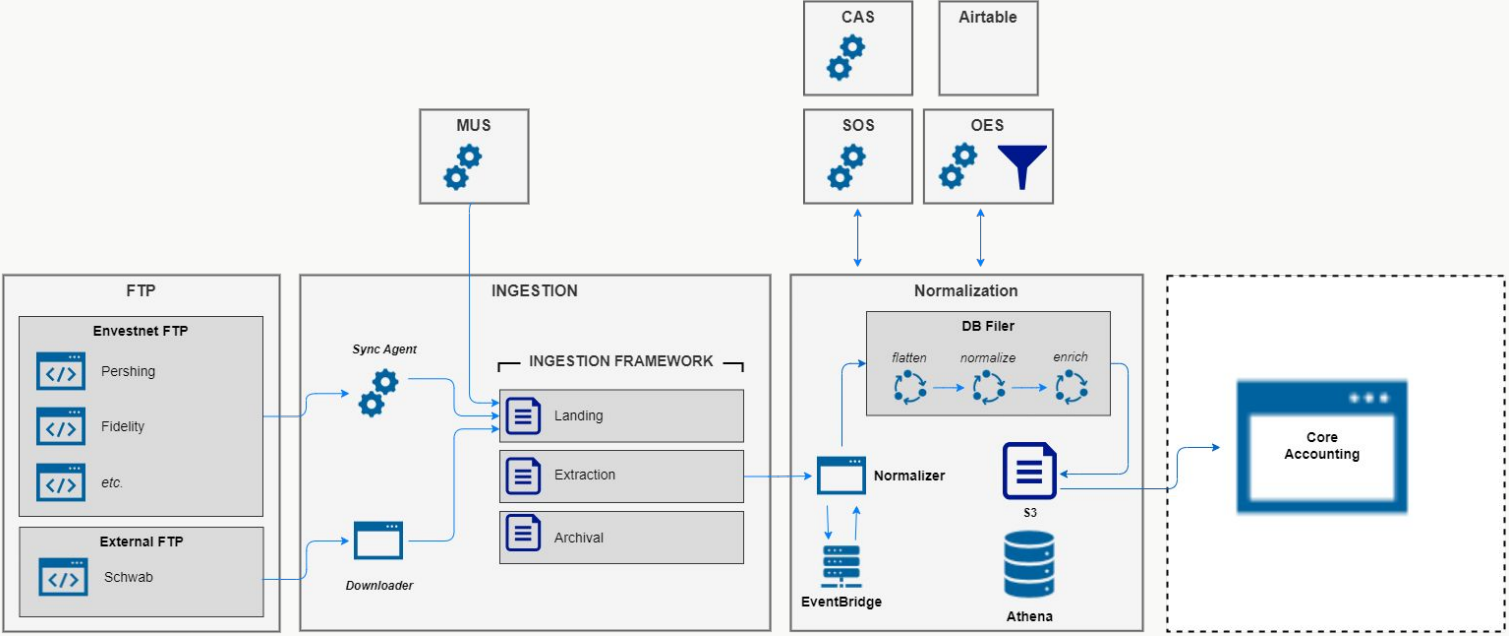
(Previously, Envestnet)



Use Case



Data Source Onboarding



Data Source Onboarding

- Onboard data from different sources
 - Batch based and event-based sources
 - Files of different formats, types, size, layout, arrival schedules
- Standardize into a common data model
 - Source specific and source agnostic business logic
 - Metadata management
 - Integration with internal systems
 - Data mapping and transformation
- Deliver data to stakeholders within SLA
 - Filter out non-specific data
 - Process and deliver within 20 minutes of arrival
 - Data validation and reconciliation

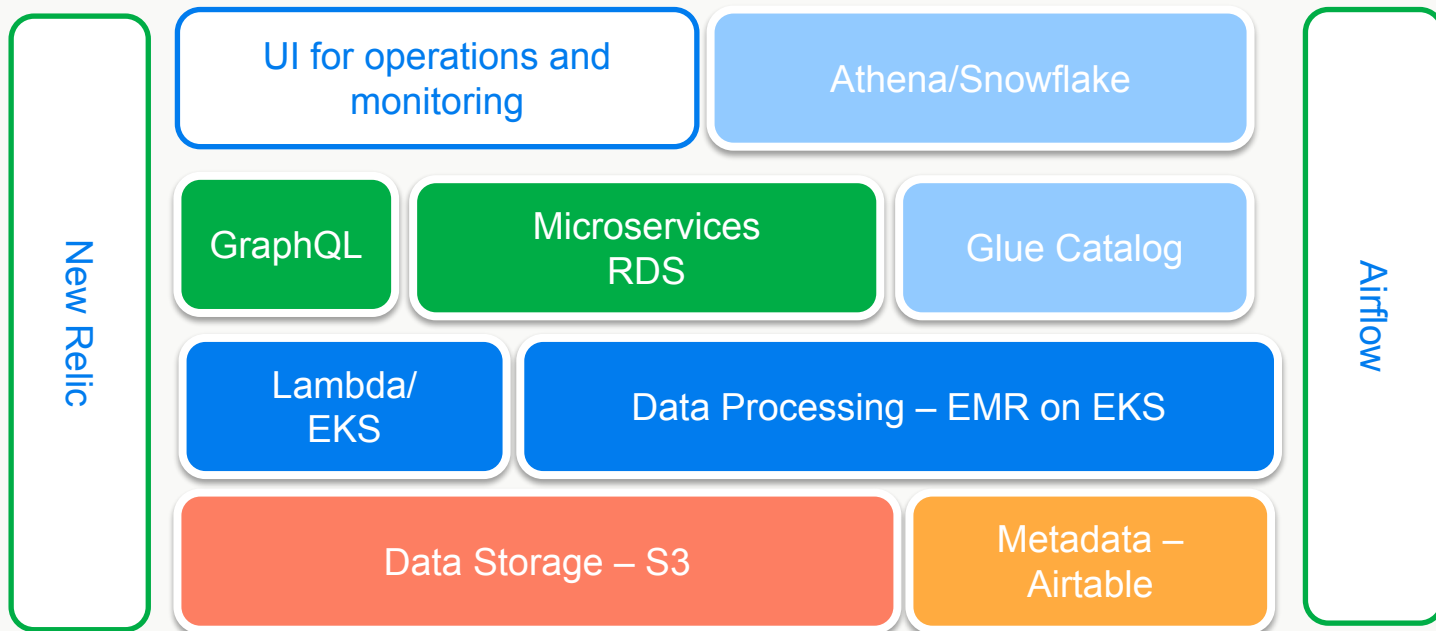


Requirements and Complexity

- **Variety:**
 - Sources have different file formats, types, layouts, sizes, file arrival patterns
- **Volume:**
 - 200 sources, 100k's of files.
 - Complete daily processing terabyte scale
- **Domain Expertise:**
 - Source data needs to be mapped and transformed to standard data model. Analysts have gained knowledge over the decade.
- **Daily Ops:**
 - Accurate and timely monitoring to run, debug & rerun before SLA
- **Processing:**
 - SLA of 20 mins. Complex multi-step processing as soon as any file dependency is met.



Technology Stack

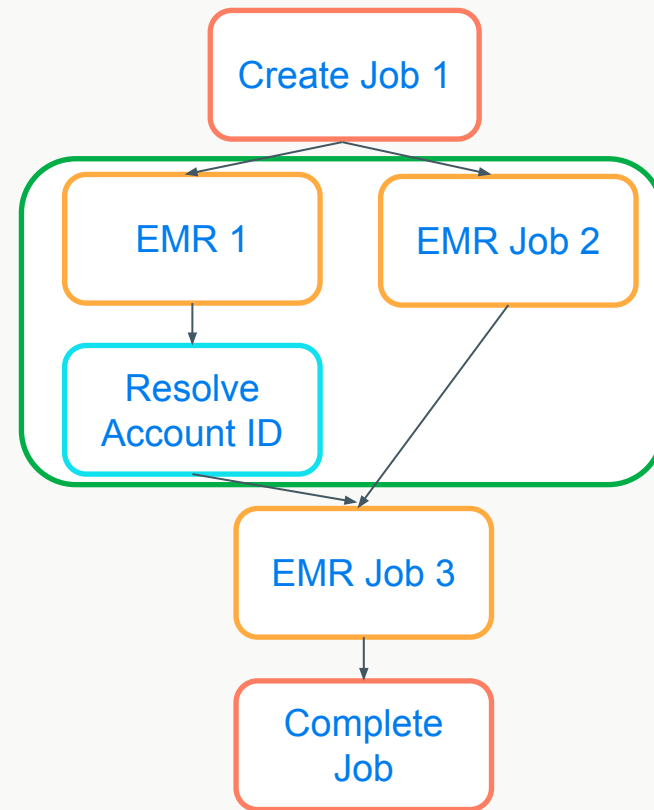
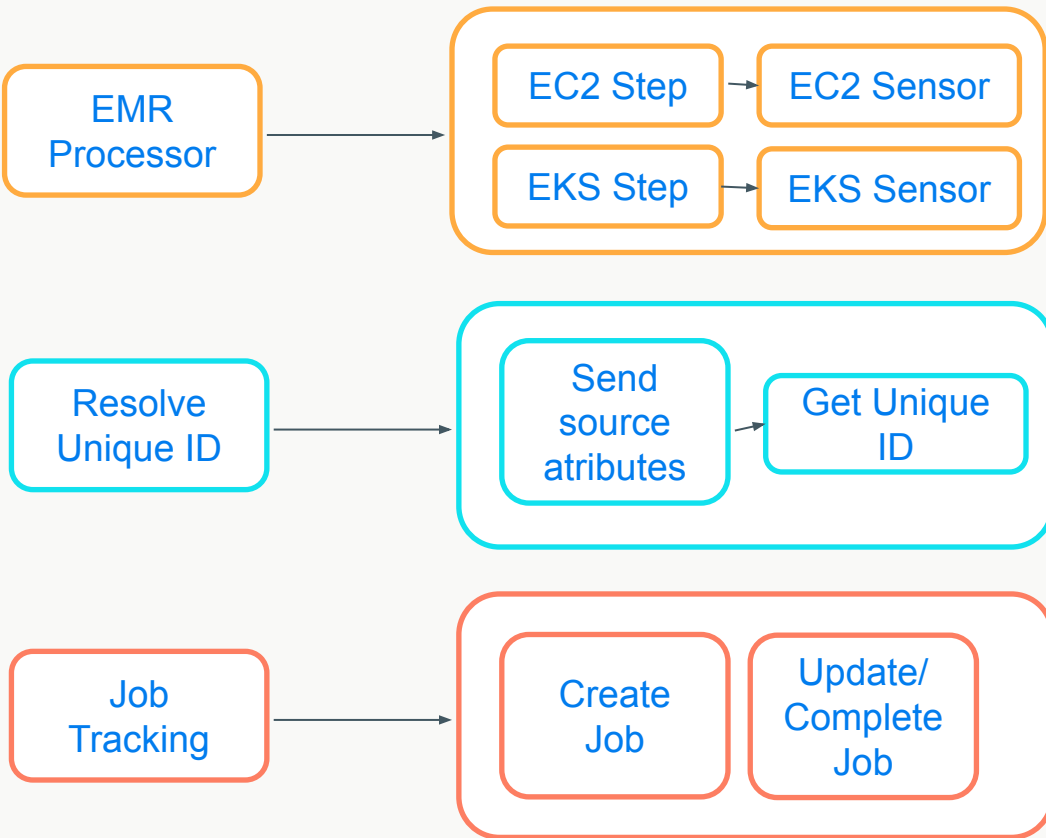
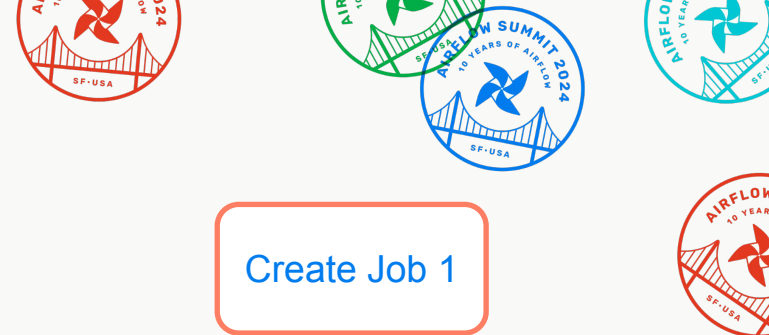


Why Airflow?

- Flexible Development
 - Scalable DAG building
 - Programmatic approach, developer friendly
 - Dynamic pipeline creation
- Easy Integration
 - Easily integrate with any system
 - Ability to develop custom functionality by extending existing functionality
- Intuitive Monitoring
 - Easy to execute, debug and monitor pipelines
 - Restarting a DAG from failure point from middle of a pipeline in case of errors
- All in One Visualization
 - Convenient dependency management between tasks
 - Visualize complex orchestration in a single view. Consolidate multiple tasks in groups.



Scalable DAG: Framework



Event Driven DAG Handling

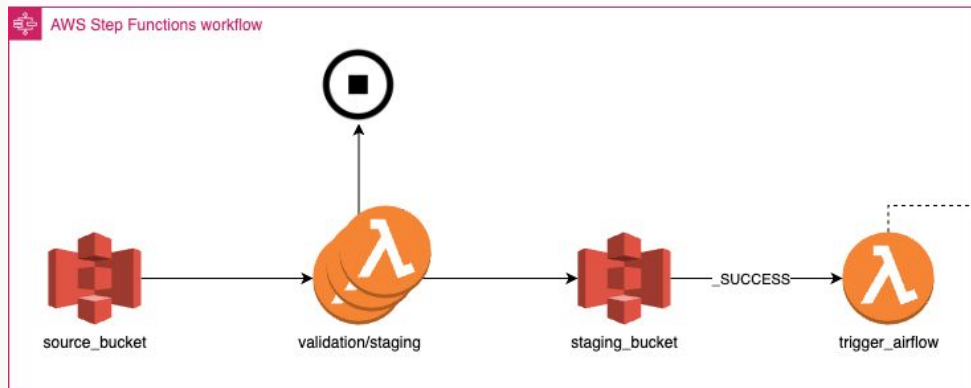


Creating Event Driven DAGs

- One DAG per source to allow monitoring per source
- Setup:
 - A configuration driven AWS Step Function that validates the files, stages them and triggers Airflow of the respective source in an event-based manner
 - Source specific configs driven by custom DAG creation modules to create DAGs automatically
 - Source and file specific & agnostic configs driven by custom spark application that contain business logic
- Airflow:
 - file_sensing task group:
 - Tasks corresponding to each file in a source
 - file_processing task group:
 - Tasks to handle file processing via API calls and spark jobs



Triggering Event Driven DAGs



`pershing_20240911_1_202409115768`

- Check if corresponding DAG execution exists (by run_id)
 - If not exists, trigger a dag execution with specific run_id
- If the execution does exist, trigger task status update via API. Mark the "task" success
- Do this step by lambda instead of sensors because sensors will constantly poll in set interval. However there is no guaranteed time arrival of files, so tasks work better



Triggering Event Driven DAGs

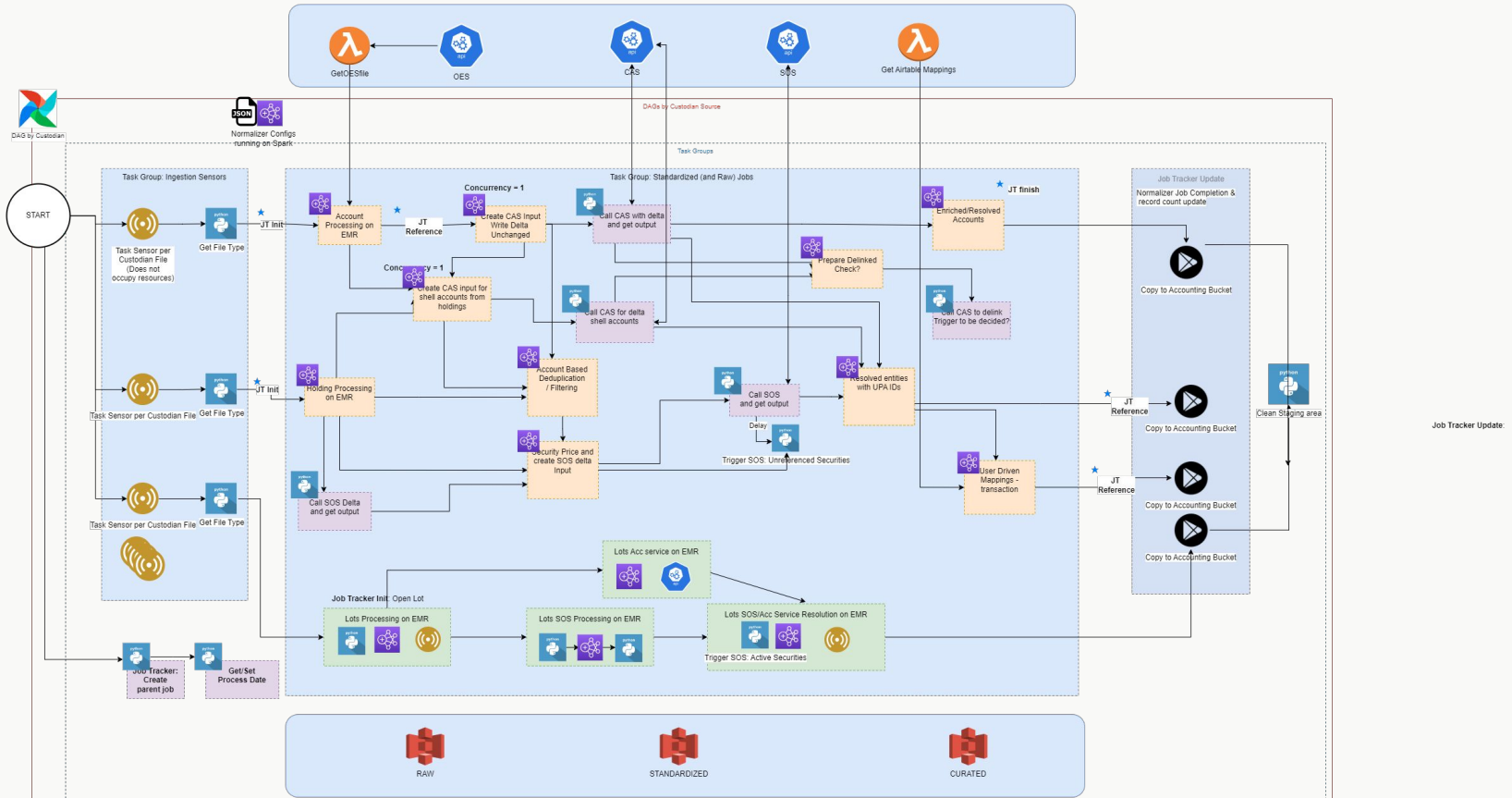
- Run IDs of the type: `<source>_<file-date>_<set_code>_<run-ts>`
- Check if an execution of that source, file_date, set_code, run_date exists or not
- If exists:
 - Mark the task for that file success via API call
- If not exists:
 - Create an execution with defined run_id of `<source>_<file-date>_<set_code>_<ts>`
 - Mark the task for that file success via API call
- All the “source” level information is passed to the DAG at conf level when triggering the DAG
- All the file level information is passed to the DAG via `_SUCCESS` file



Scalable Development



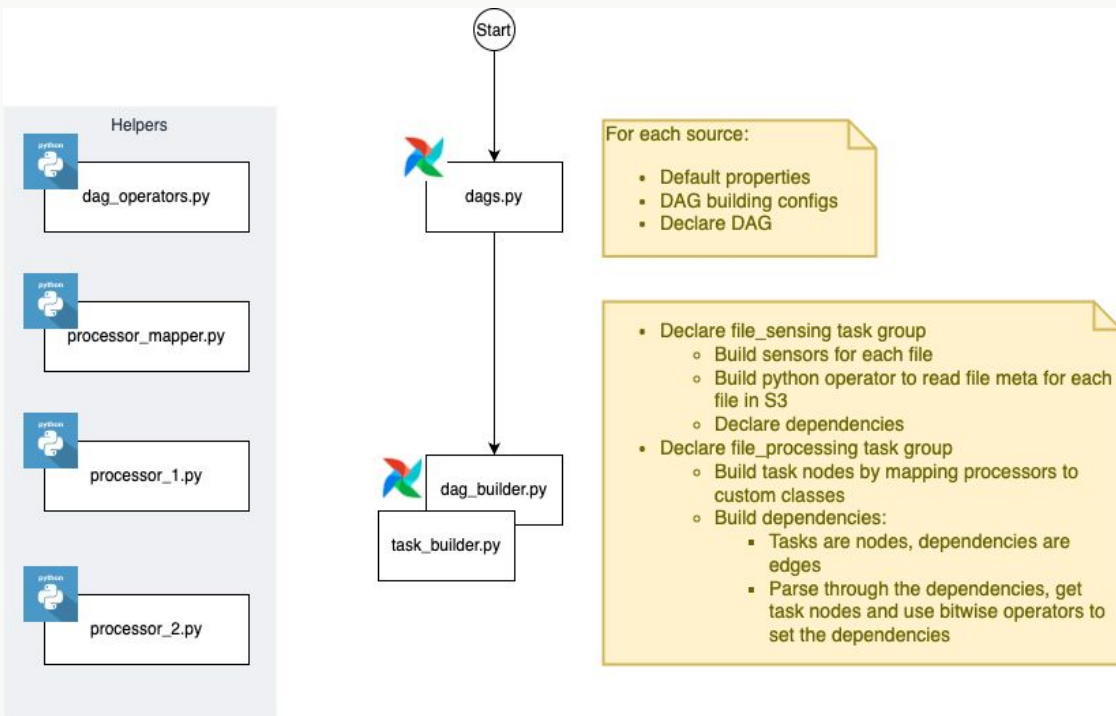
Simplified Architecture Diagram



Scalable DAG: Config Building

Dependency Management Config

```
"Call API 1": {  
  "processor": "call_api_1",  
  "label": "Call API to get output in S3 bucket",  
  "task_group": "file_processing",  
  "dependencies": {  
    "file_sensors": [  
      "file1",  
      "file2"  
    ]  
  }  
},  
"Submit Spark Job": {  
  "processor": "submit_spark_job",  
  "label": "Submit Spark Job for <source>",  
  "task_group": "file_processing",  
  "dependencies": {  
    "file_processing": [  
      "call_api_1"  
    ]  
  }  
}
```



Scalable DAG: Config Building



Dependency Management Config

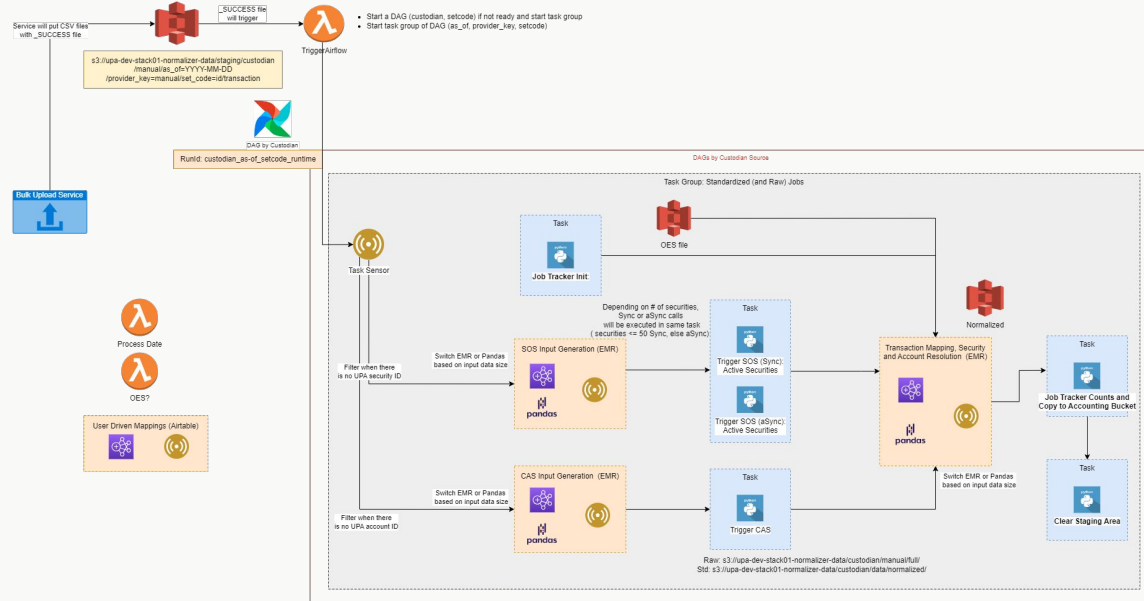
```
"call_api_1": {
  "processor": "API",
  "label": "Call API to get output in S3 bucket",
  "task_group": "file_processing",
  "dependencies": {
    "file_sensors": [
      "file1",
      "file2"
    ]
  }
},

"submit_spark_job": {
  "processor": "EMR",
  "label": "Submit Spark Job for <source>",
  "task_group": "file_processing",
  "dependencies": {
    "file_processing": [
      "call_api_1"
    ]
  }
}
```

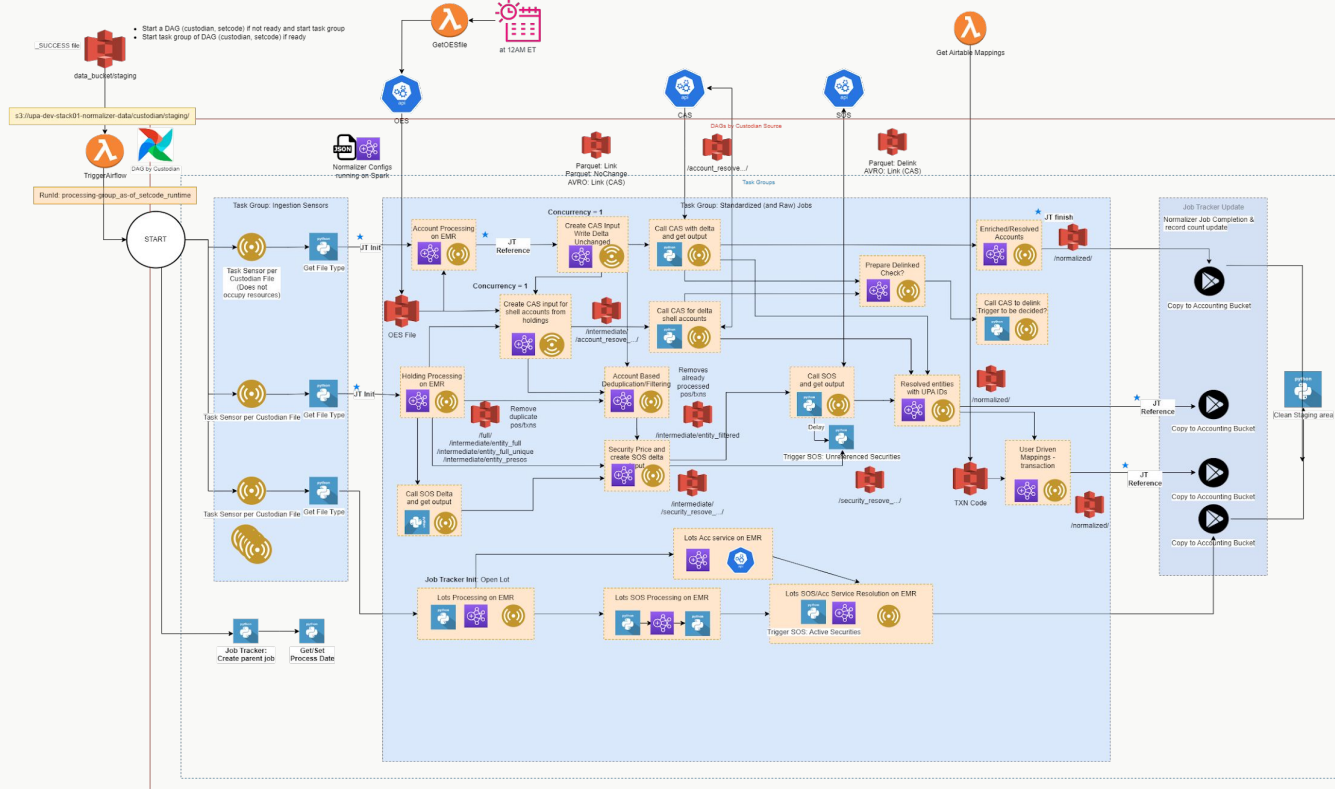
```
for parent in dependencies:
    first_task = task.steps[0]
    if stage == 'root':
        dag.get_task(parent) >> Label(edge_label) >> first_task
    elif stage == destination_stage:
        tasks[parent].steps[-1] >> Label(edge_label) >> first_task
    else:
        if dag.has_task(f"{parent_task_group.group_id}.{parent}{parent_task_suffix}"):
            parent_task_group.get_child_by_label(
                f"{parent}{parent_task_suffix}") >> Label(edge_label) >> first_task
```


Simple DAG

MUS-Normalizer Pipeline Architecture



Complex DAG



Questions?

Email:

ipsa.trivedi@gmail.com

connect.subramanian@gmail.com

LinkedIn:

<https://www.linkedin.com/in/ipsatrivedi>

<https://www.linkedin.com/in/svellaiyan>

