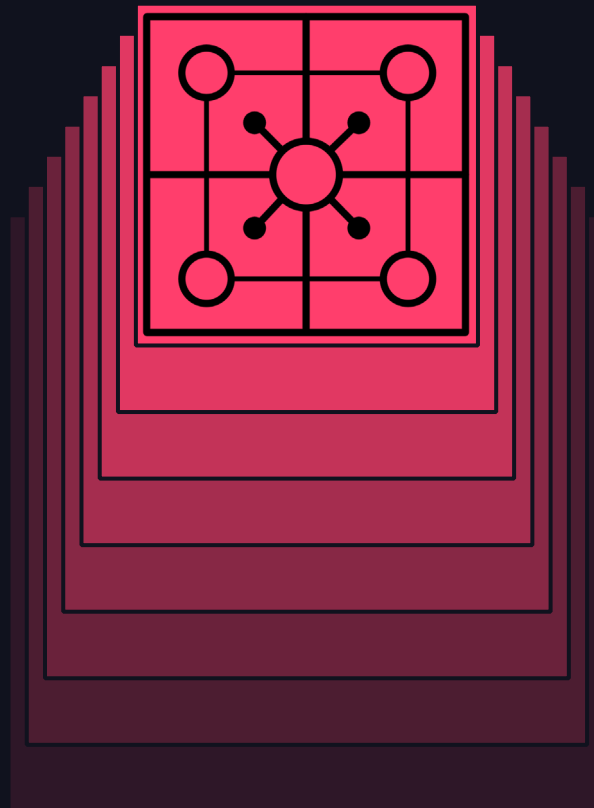


BLUEPRINT



Roberto Alejandro Flores Meregote
2024-06-08

NICE TO MEET YOU



BLUEPRINT

THE DE TEAM!



BLUEPRINT

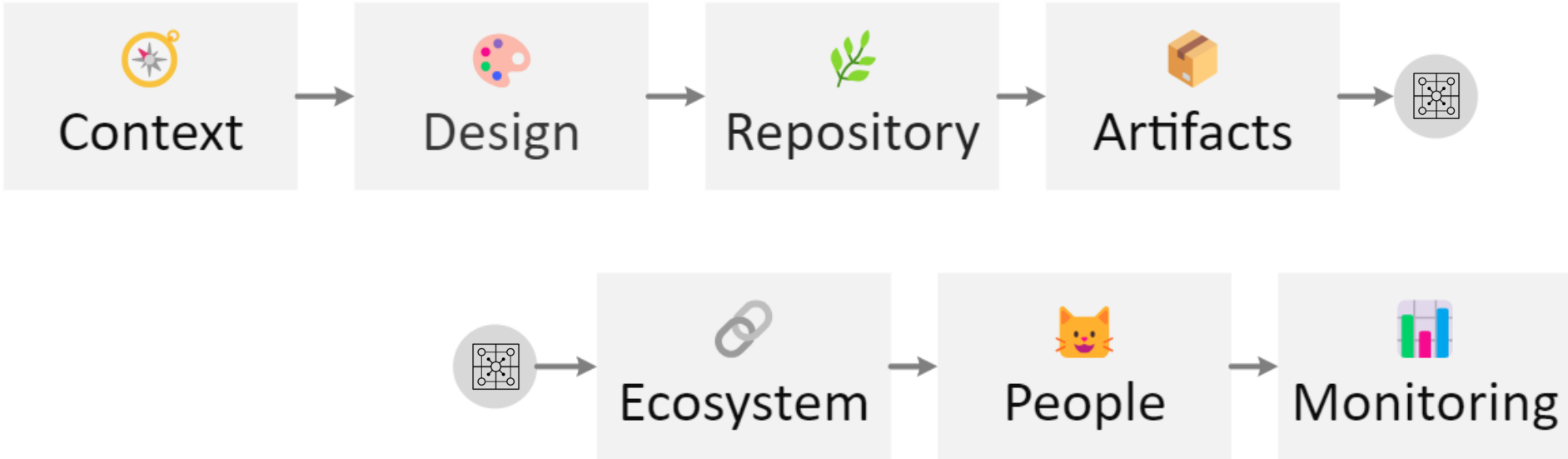


Roberto Flores Meregote
Europe Head of Data Engineering
@Unilever



BLUEPRINT

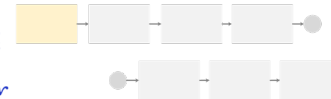
AGENDA!



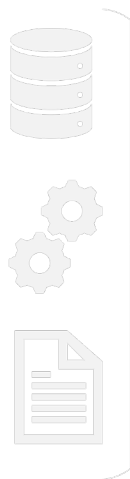
CONTEXT



UNILEVER'S DATA ESTATE



*Internal & External
Data Sources*



UDL (Universal
Data Lake) -
copy of source
in full fidelity
where all
enterprise data
sources land

BDL (Business
Data Lake) -
segmented by
function
(marketing,
finance, supply
chain) with
applied
business rules

MDL (Market
Data Lake) -
built to service
region needs
across
functions

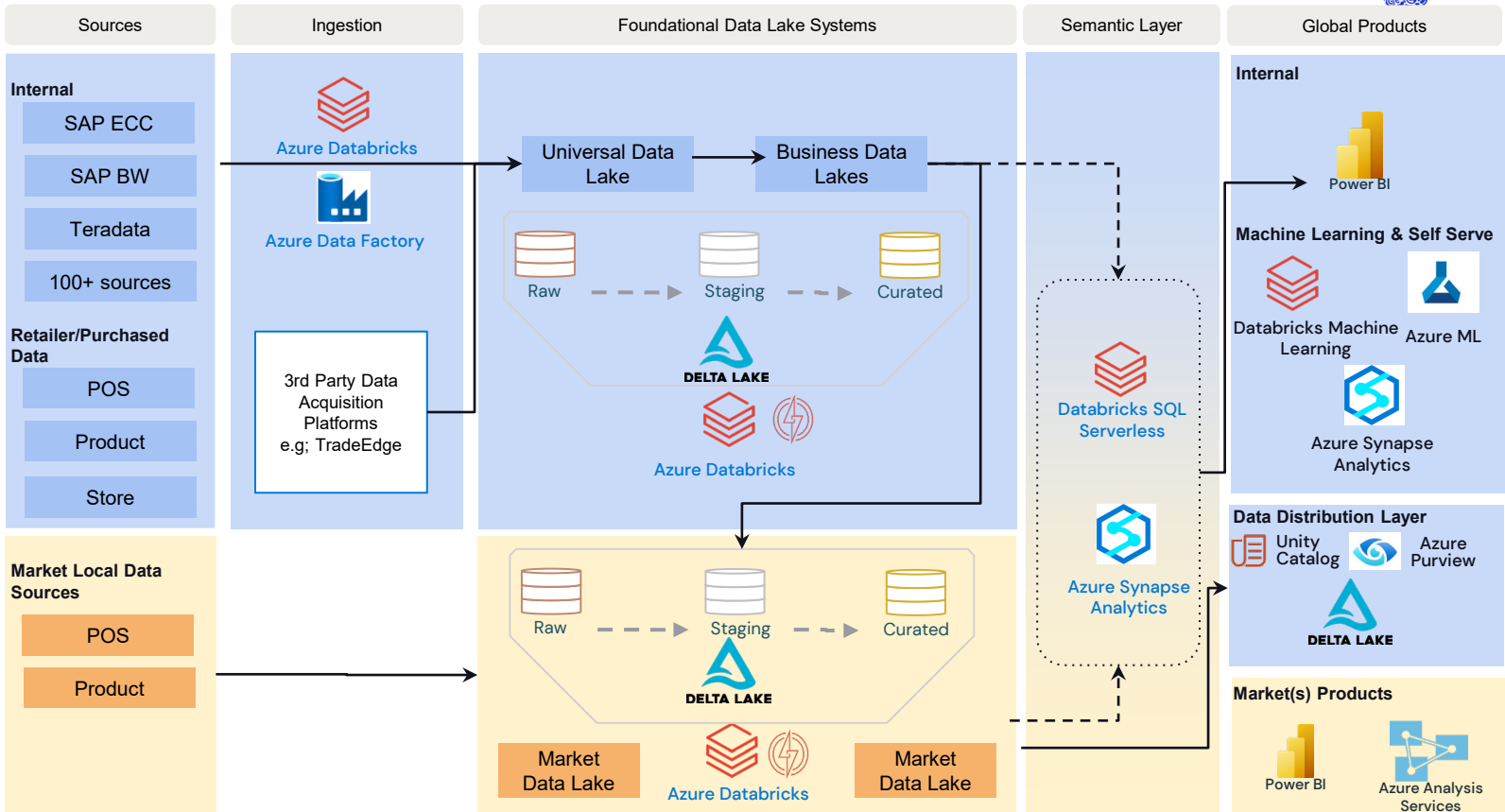
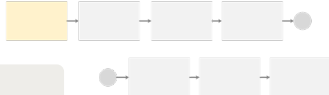
PDS (Product
Data Store) -
product or
application
specific logic
built usually for
a smaller
audience and
focused
purpose



End Users



UNILEVER'S DATA ESTATE



Unity Catalog



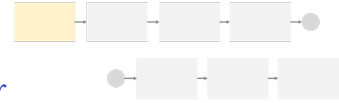
Azure Active Directory



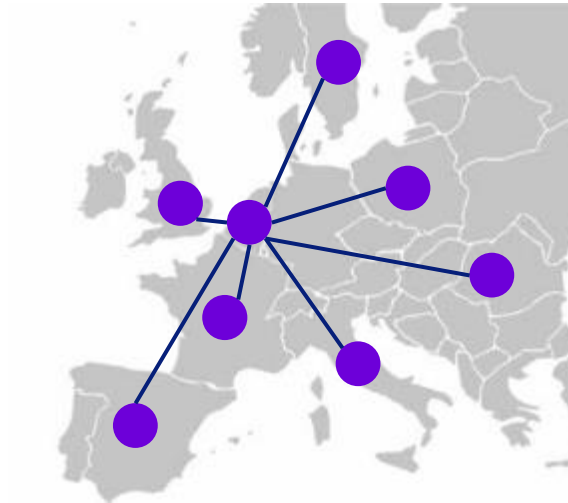
Azure Purview



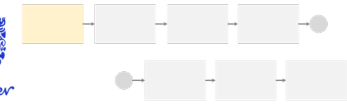
EUROPE IS A COMPLEX MARKET





PRIOR TO MDL, COUNTRIES HAD DIFFERENT DATA MATURITY LEVELS



EUROPE IS A COMPLEX MARKET



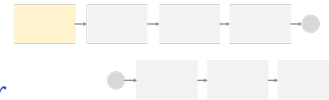
SIMILAR IN SIZE TO NORTH AMERICA, BUT WITH HIGH DATA COMPLEXITY

| DATA COMPLEXITY | | |
|----------------------|---|---|
| |  |  |
| UL Markets | 2 | 38 |
| UL BG Cells | 10 | 129 |
| Official Languages | 1 | 24 |
| Databases (external) | 20 | 316 |

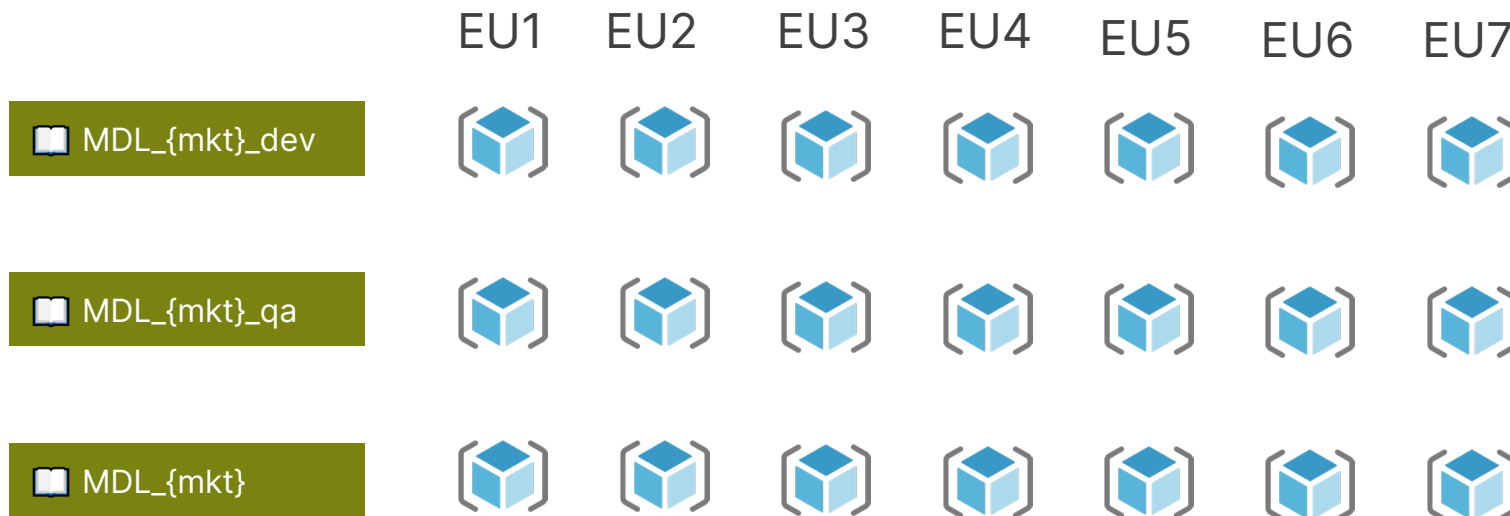


Source: Europe MDL | Country / Category Cells with a Nielsen subscription

EUROPE IS A COMPLEX MARKET



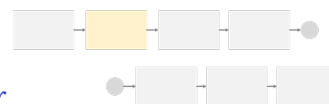
PREVIOUS SETUP MEANT THAT WE STARTED WITH MANY ENVIRONMENTS



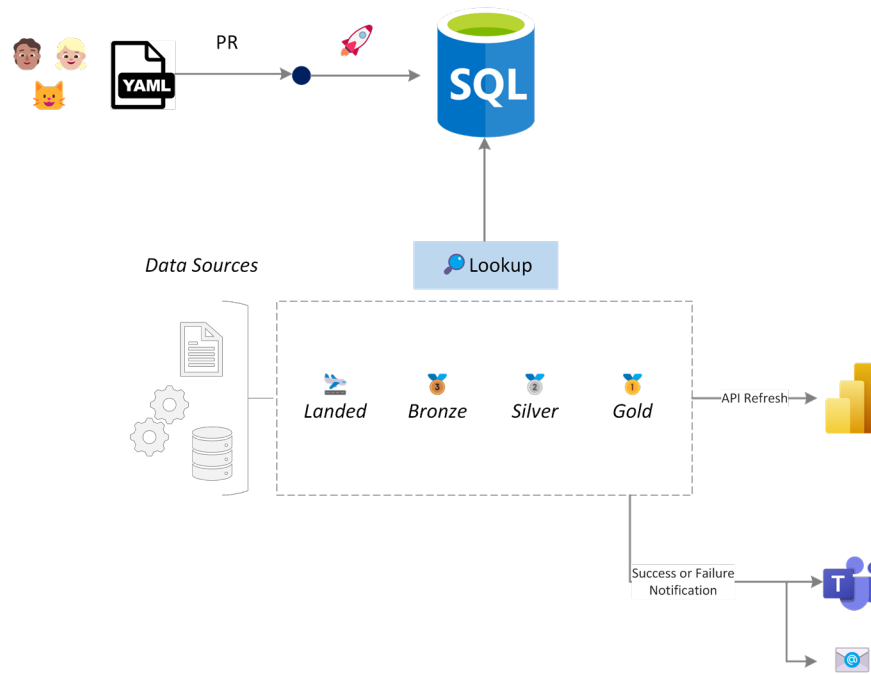
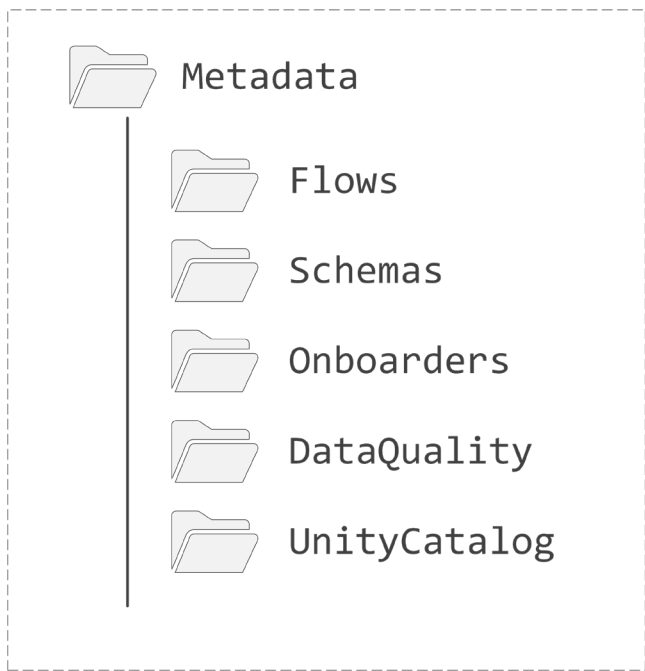
FRAMEWORK



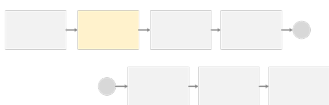
FRAMEWORK OVERVIEW



WHAT WE WANTED TO MANAGE IN A TECH AGNOSTIC WAY



FRAMEWORK OVERVIEW



WHAT WE WANTED TO MANAGE IN A TECH AGNOSTIC WAY

Common Artifacts

1. Onboarding
2. Unity Catalog Management
3. Power Platform Integrations
4. Data Factory Pipelines

← Blueprint 2.2.6


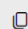
Overview Versions

↑ Promote

× Delete

⚙ Allow externally-sourced versions

Get this package

 Connect to Feed then `pip install Blueprint==2.2.6` 


Summary

Library for managing and processing data assets. Version above 2.0.11 furnance is not backward compatible due to parameter changes

Provenance

Build

 [903446-blueprint-build-whl / 20240521.1](#)

 May 21 (5/21/2024 2:26:31 PM)


Code

 [af72158e](#) in  [main](#)

 [903446-blueprint](#)

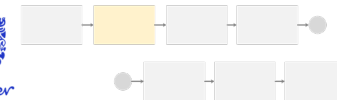
Development

Source

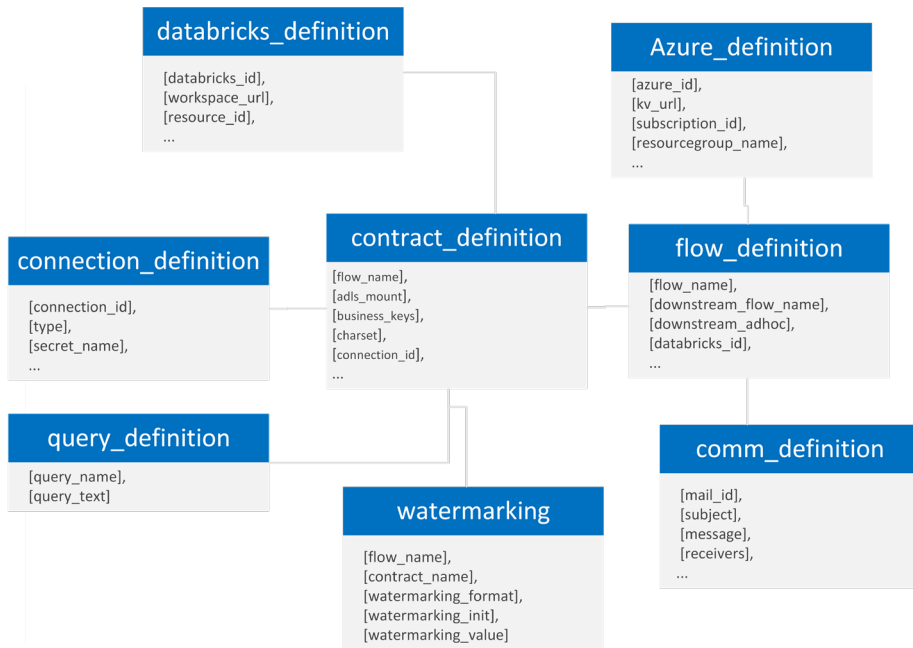
 This feed



FRAMEWORK OVERVIEW



DATABASE

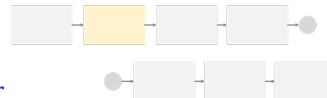


Stored Procedure Sample

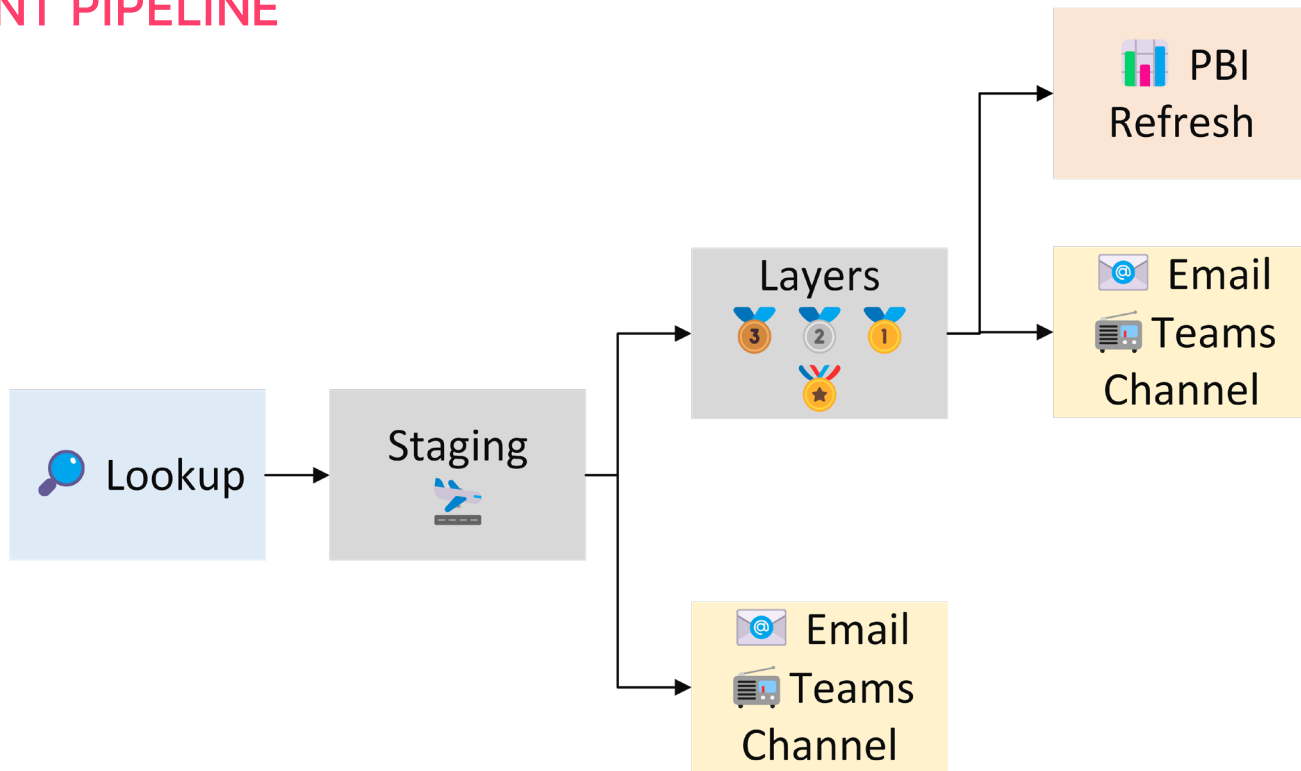
```
CREATE OR ALTER PROCEDURE [utils].[sp_staging_info]
    @flow_group_name NVARCHAR(255),
    @staging_type NVARCHAR(255),
    @meta_schema NVARCHAR(255),
    @flow_filter NVARCHAR(255) = NULL
AS
SELECT DISTINCT
fw.flow_name,
fw.blueprint_version,
az.kv_url,
az.landing_adls_url,
az.delta_lake_path,
az.subscription_id,
az.tenant_id,
az.resourcegroup_name,
az.secret_scope,
az.azure_devops_organization,
ct.user_id,
ct.user_name,
(...)
```



WHAT DOES IT LOOK LIKE

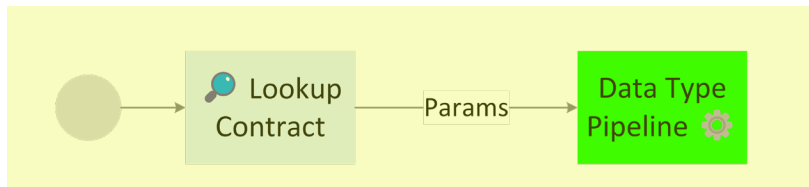
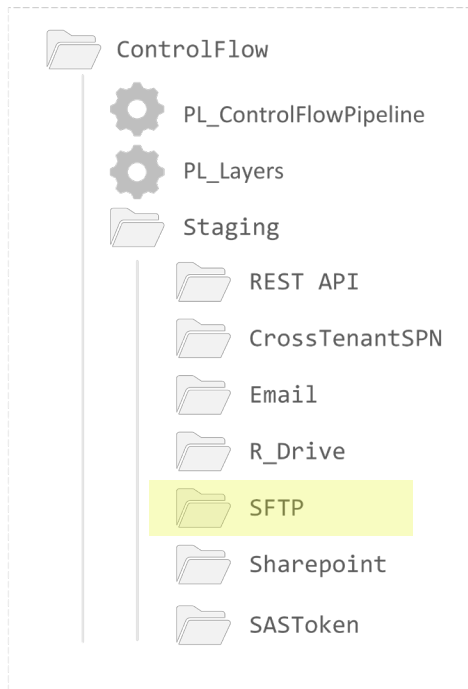


PARENT PIPELINE



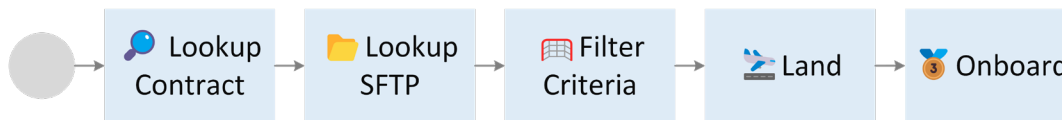
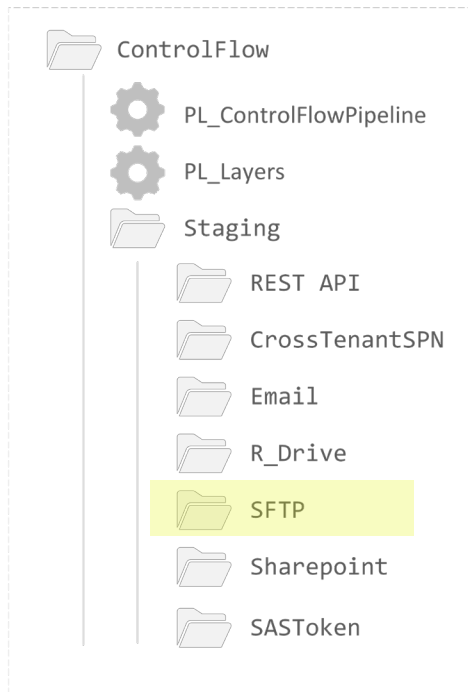
WHAT DOES IT LOOK LIKE

STAGING PIPELINES

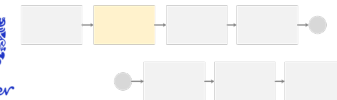


WHAT DOES IT LOOK LIKE

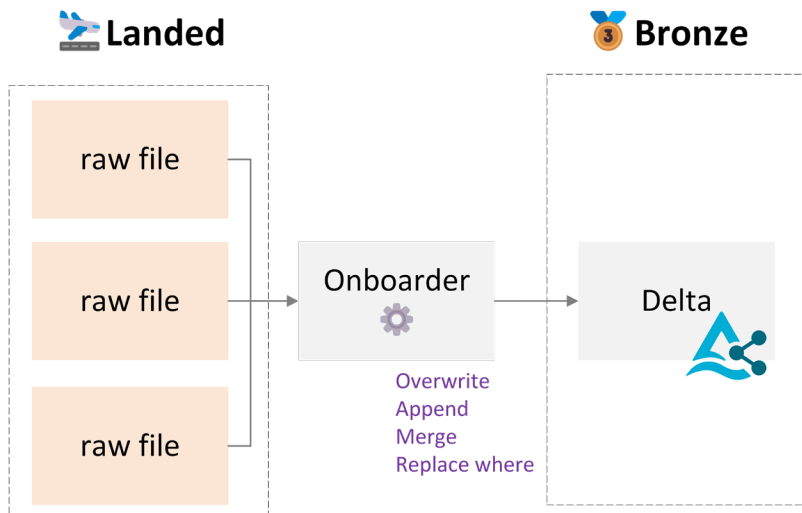
STAGING PIPELINE SAMPLE - SFTP



WHAT DOES IT LOOK LIKE

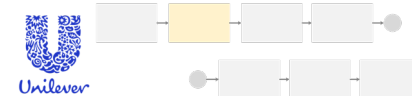


ONBOARDING



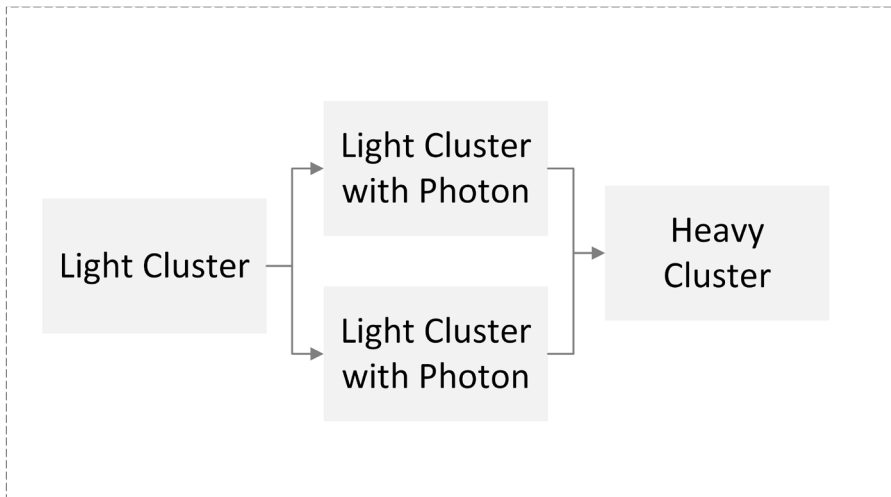
```
- contract_name: BNL_PDD_LOOKUP_LIST_INGRED
adls_mount: /dbfs/mnt/landingzone/
business_keys: null
connection_id: SAP_DQ_SP
delimiter: null
file_type: .xlsx
header: 1
layer: bronze
load_type: F
mapping_id: null
onFail: Fail
range: A1
region: Benelux
schema_name: bronze_md1_bn1sharepoint
sheet_name: ingredientGPC
skip_line: null
source_file_name: lookup_lists.xlsx
source_path: sites/PDDsourcedata/Shared%20Documents/PDD%20lookup%20lists
source_system_type: AdHoc
source_system: Sharepoint
stagingType: PL_Sharepoint_MetadataBased
table_name: bdl_pdd_lov_ingred
```

WHAT DOES IT LOOK LIKE



LAYER EXECUTION, DEPENDENCIES AND RUN ORDER

 Gold



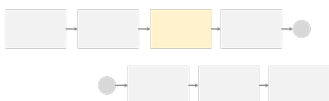
```
gold_notebooks:
| #Finance Dim
- path: /MarketDataLake/PRO/05-Gold/Finance/gold_finance_dim
  flow_filter: Finance
  parameters:
  | DefaultParm: NA
  run_order: 1
  databricks_id: Normal Cluster
| #Nielsen Dim
- path: /MarketDataLake/PRO/05-Gold/Nielsen/gold_nielsen_dim
  flow_filter: Poland
  parameters:
  | param_country_2: Poland
  | param_country_cd: PL
  databricks_id: Heavy Cluster
  run_order: 2
- path: /MarketDataLake/PRO/05-Gold/Nielsen/gold_nielsen_dim
  flow_filter: Hungary
  parameters:
  | param_country_2: Hungary
  | param_country_cd: HU
  run_order: 2
  databricks_id: Normal Cluster
```



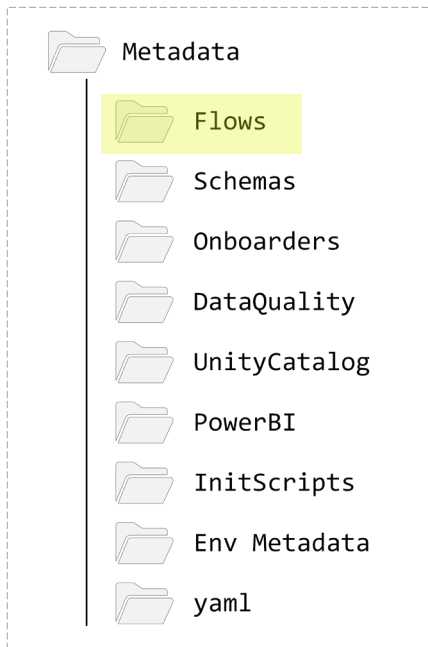
REPO STRUCTURE



METADATA REPOSITORY STRUCTURE



FLOW FILES

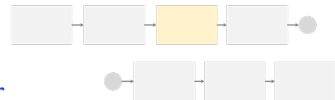


sample: Flow.yaml

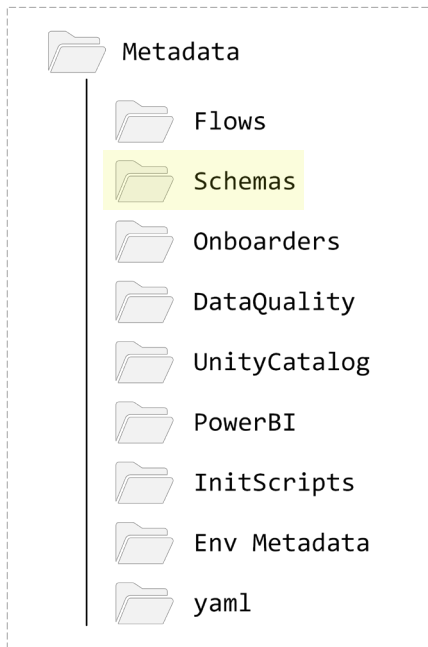
```
Flows:
- flow_name: DACH_Retailer
  azure_id: Azure Info
  bi_id: 'EPOS_DE_PBI'
  databricks_id: UC_Enabled_Cluster
  downstream_adhoc: ''
  downstream_flow_name: ''
  mail_id: 'DACH ePos Mail'
  run_downstream: True
  master_pipeline: PL_Master_MetadataDriven
  blueprint_version: 2.0.12
Contracts:
- contract_name: DACH_REWE_SALES_AS
  adls_mount: /dbfs/mnt/landingzone/
  business_keys: null
  connection_id: DACH SFTP
  delimiter: ;
  file_type: .csv
  header: 0
  layer: bronze
  (...)
```



METADATA REPOSITORY STRUCTURE



SCHEMAS

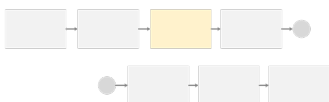


Sample: Schemas.json

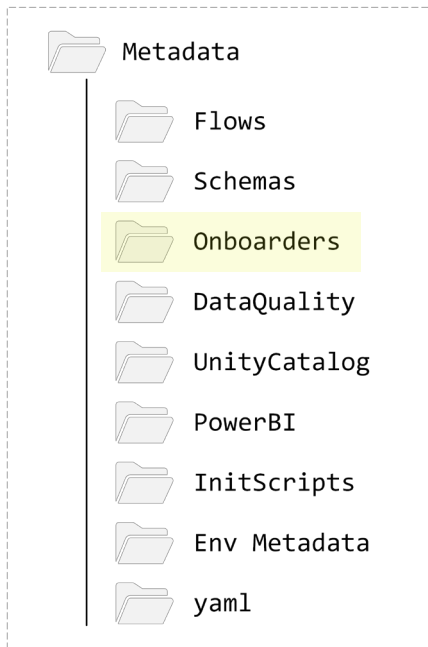
```
{
  "$schema": "http://json-schema.org/draft-07/schema#",
  "type": "object",
  "properties": {
    "Flows": {
      "type": "array",
      "items": {
        "type": "object",
        "properties": {
          "flow_name": {
            "type": ["integer", "string"]
          },
          "downstream_flow_name": {
            "type": ["string", "null"]
          },
          "downstream_adhoc": {
            "type": ["string", "null"]
          }
        }
      }
    }
  }
  (...)
}
```



METADATA REPOSITORY STRUCTURE



GENERALIZED ONBOARDER



Sample: Onboarder.py

```
# Databricks notebook source
from blueprint.engineering.kiln import Kiln
import os
# COMMAND -----
dbutils.widgets.removeAll()

flow_name = dbutils.widgets.get("flow_name")
contract_name = dbutils.widgets.get("contract_name")
meta_schema = dbutils.widgets.get("meta_schema")

print(f"Flow Name: {flow_name}")
print(f"Onboarding Contract: {contract_name}")
print(f"Metadata Schema: {meta_schema}")

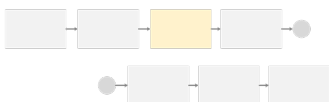
# COMMAND -----

print(f'Running with version of Blueprint:
{os.environ["PKG_VERSION"]}')

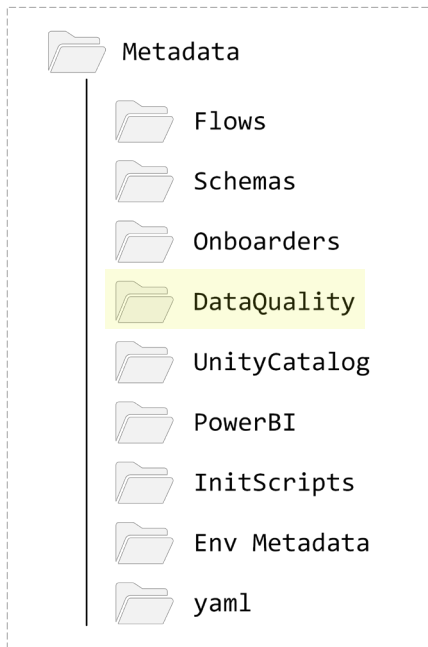
(...)
```



METADATA REPOSITORY STRUCTURE



DATA QUALITY (GREAT EXPECTATIONS)



Sample: DataQuality.py

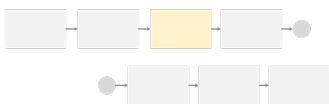
```
import pyspark.sql.functions as F
import json
from pyspark.sql.types import IntegerType, StringType, DoubleType, LongType
from delta.tables import *
from great_expectations.core.batch import RuntimeBatchRequest
from great_expectations.util import get_context
from great_expectations.data_context.types.base import (
    DataContextConfig,
    FilesystemStoreBackendDefaults,
)

import great_expectations as ge
from great_expectations.dataset import (
    SparkDFDataset,
    MetaSparkDFDataset,
)

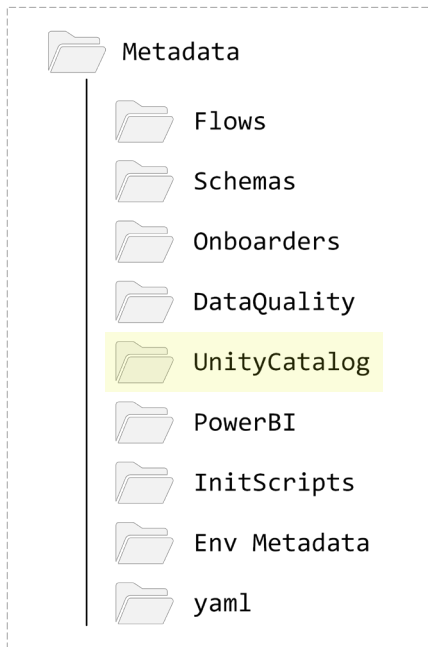
(...)
```



METADATA REPOSITORY STRUCTURE



UNITY CATALOG METADATA

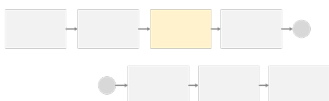


Sample: UnityCatalogDefinition.yaml

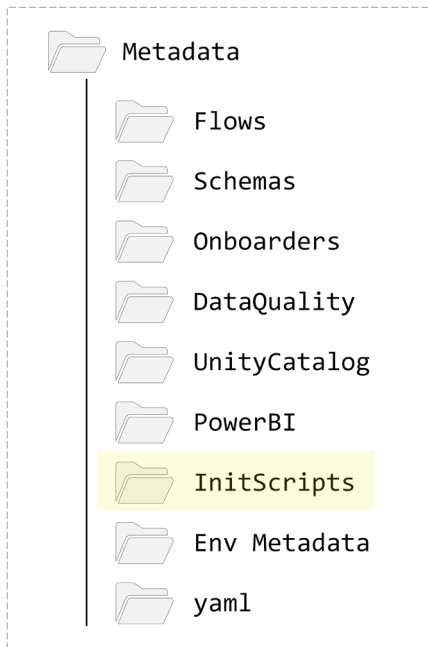
```
UC_Def:
- flow_name: Init_schema
  schemas:
  - name: init_schema
    comment: "This is just a dummy file to init the correct schema"
    tags:
    - Test
    access:
    - acls_group_name: "SEC-ES-DA-p-903444-europe-analyst"
      acls_access_type: "SELECT, MODIFY"
      acls_access_catalogue: "mdl_europe_anz_dev"
    tables:
    - name: init_she
      comment: "### Supports Markdown 1. First item 2. Second item"
      tags:
      - Test          columns:
        - name: "Country_Code"
          comment: "Country code key as of 07/03/2024"
          tags:
          - PK
```



METADATA REPOSITORY STRUCTURE



INIT SCRIPTS



Sample: BlueprintInit.sh

```
curl -sL https://aka.ms/InstallAzureCLIDeb | bash

python -m pip install --upgrade pip

az login --service-principal -u ${AZ_DEVOPS_SP_APP_ID} -p
${AZ_DEVOPS_SP_SECRET} --tenant ${AZ_TENANT_ID}

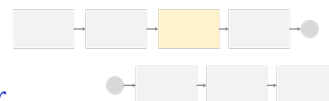
AZ_SP_TOKEN=$(az account get-access-token --query accessToken -o tsv)

pip install Blueprint==${PKG_VERSION} --index-url
"https://token:${AZ_SP_TOKEN}@pkgs.dev.azure.com/${AZ_DEVOPS_ORG_NAME}/_packaging/${AZ_DEVOPS_FEED_NAME}/pypi/simple/"

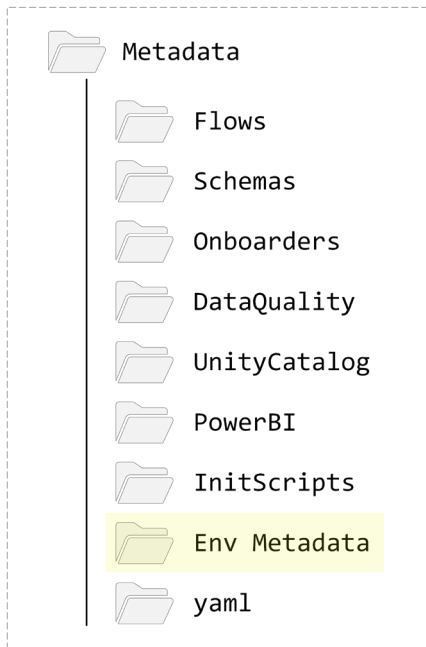
(...)
```



METADATA REPOSITORY STRUCTURE



UNITY CATALOG METADATA



Sample: EnvMetadata_dev.yaml

```
DatabricksDefinition:
- databricks_id: Strong Cluster
  default_catalog: mdl_europe_anz_dev
  workspace_url: https://adb-xxxxxxx.yy.azuredatabricks.net
  resource_id: /subscriptions/blablabla
  policy_id: 0000EB555555
  cluster_version: 12.2.x-scala2.12
  cluster_type: Standard_E8d_v4
  scaling: '5:20'
  init_script: /Metadata/InitScripts/BlueprintInit.sh
AzureDefinition:
- azure_id: Azure Info
  kv_url: https://xxxxx.vault.azure.net/
  landing_adls_url: https://xxxxxxx.dfs.core.windows.net/
  delta_lake_path: 'abfss://unilever@
    xxxxxx.dfs.core.windows.net/datalake'
  subscription_id: bbbbb-aaaa-zzzz-yyyy-xxxxxx
  resourcegroup_name: abceasyasonetwothree-rg
  tenant_id: aaaaaa-bbbbb-ccccc-dddd
  azure_devops_organization: bnxxxxxxxxts
  (...)
```

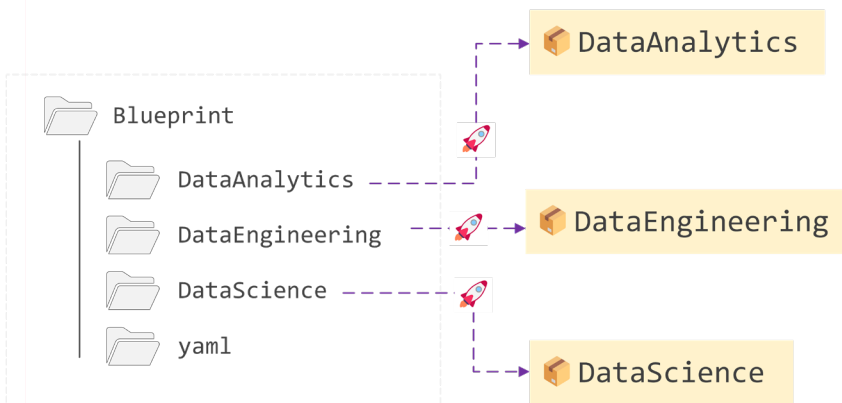


ARTIFACTS



HOW DO WE DEPLOY IT

PYTHON ARTIFACTS



YAML sample: Blueprint Python Artifact Creation

```
stages:
- stage: Building_new_version
  jobs:
  - job: Build
    pool:
      name: Azure Pipelines
      vmImage: windows-2019

steps:
- checkout: self
  persistCredentials: true
  clean: true

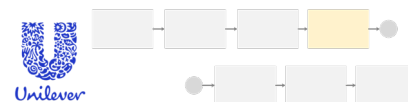
  - script: |
    pip install twine
    pip install build
    displayName: 'Install twine and build'

  - script: |
    python -m build -w blueprint
    displayName: 'Build wheel'

  - task: TwineAuthenticate@1
inputs:
  artifactFeed: Europe_MDL

  - script: |
    python -m twine upload -r Europe_MDL --config-file $(PYPIRC_PATH) blueprint/dist/*.whl
    displayName: 'Upload build whl to feed'
```

PYTHON PACKAGE ARTIFACT FEED



PACKAGE



[BluePrint](#)

Version 2.2.6

| <input type="checkbox"/> | Version | Views | Publish date | Source |
|--------------------------|--|-------|--------------|-----------|
| <input type="checkbox"/> | 2.2.6 (currently selected) | Local | Tuesday | This feed |
| <input type="checkbox"/> | 2.2.4 | Local | Apr 15 | This feed |
| <input type="checkbox"/> | 2.2.3 | Local | Apr 10 | This feed |
| <input type="checkbox"/> | 2.2.1 | Local | Apr 2 | This feed |

Get this package



Connect to Feed

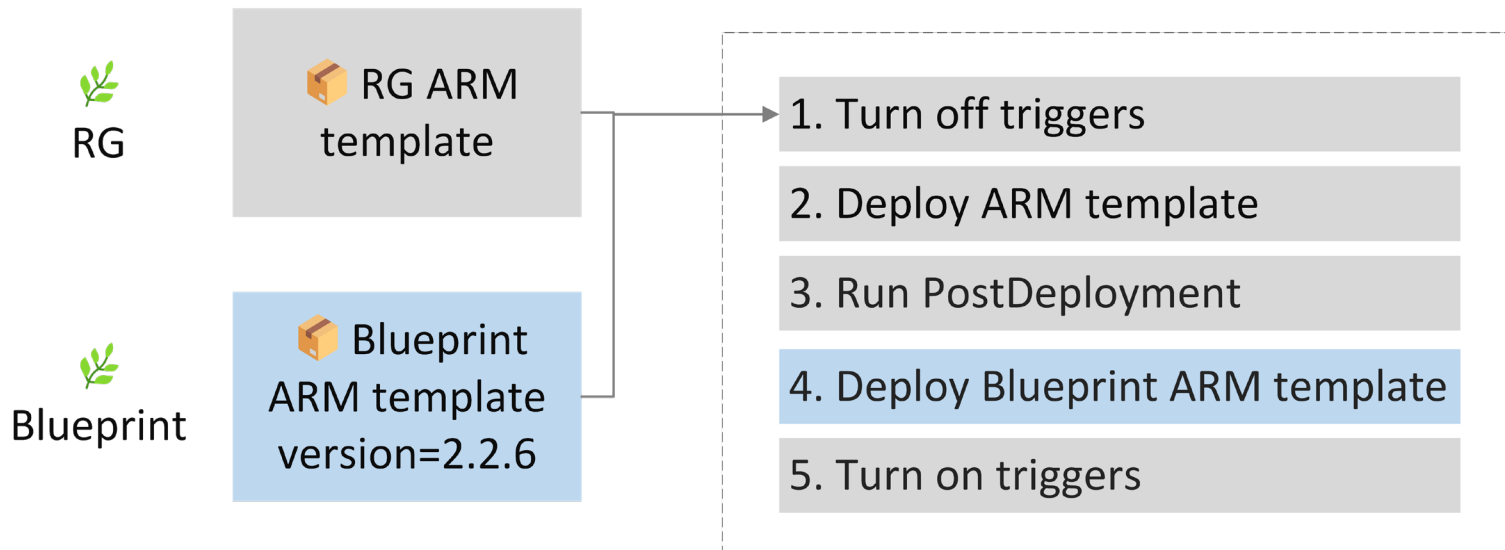
then

```
pip install BluePrint==2.0.12
```

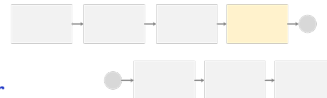


PYTHON PACKAGE ARTIFACT FEED

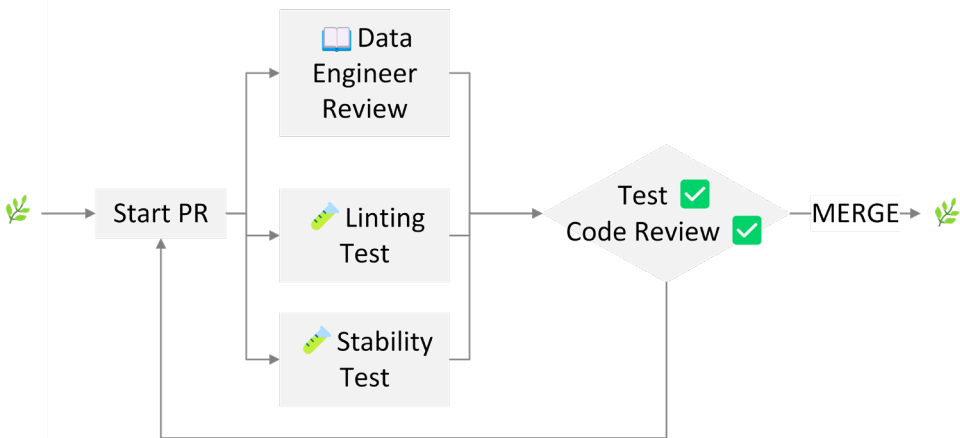
DATA FACTORY ARM TEMPLATES



REUSABILITY REQUIRES RESILIENCE



TESTING AIMS TO ENSURE STABILITY ACROSS VERSIONS 



Testing Pipeline Screenshot

Updated furnace.py Approve Set auto-complete

Active !198067 **PR** Roberto Flores proposes to merge [roberto/test-pr](#) into [main](#)

Overview Files Updates Commits Conflicts

- 1 required check running now
1 optional check not yet run
- blueprint-testvalidation blueprint-testvalidation in progress
[View 2 checks](#)
- At least 1 reviewer must approve
- No merge conflicts
Last checked Just now

Reviewers Add

Required
No required reviewers

Optional
No optional reviewers

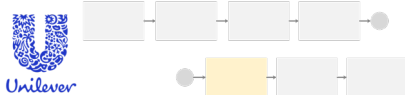
Tags +

No tags



INTEGRATION WITH DATA ESTATE

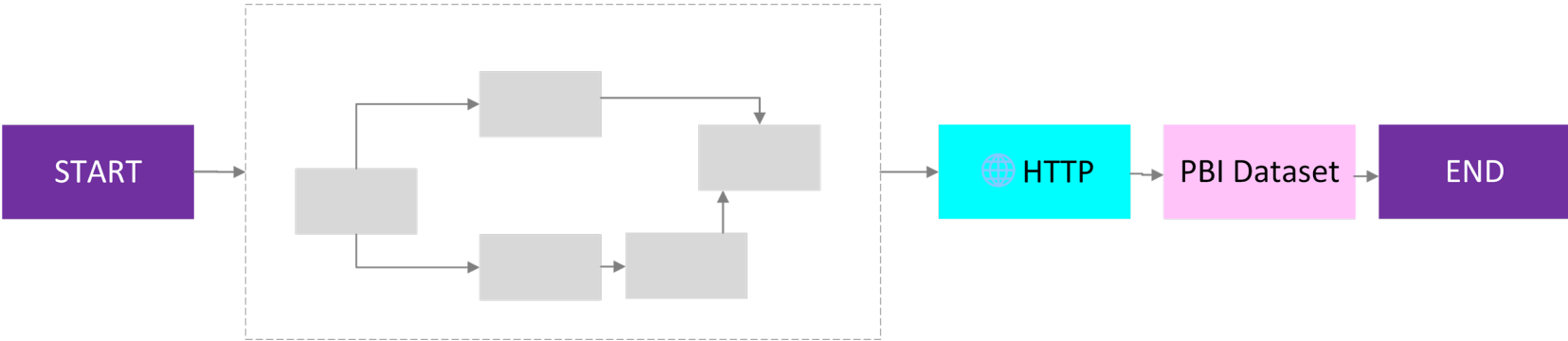
PBI PREMIUM CAPACITY



REFRESHING VIA ENHANCED API

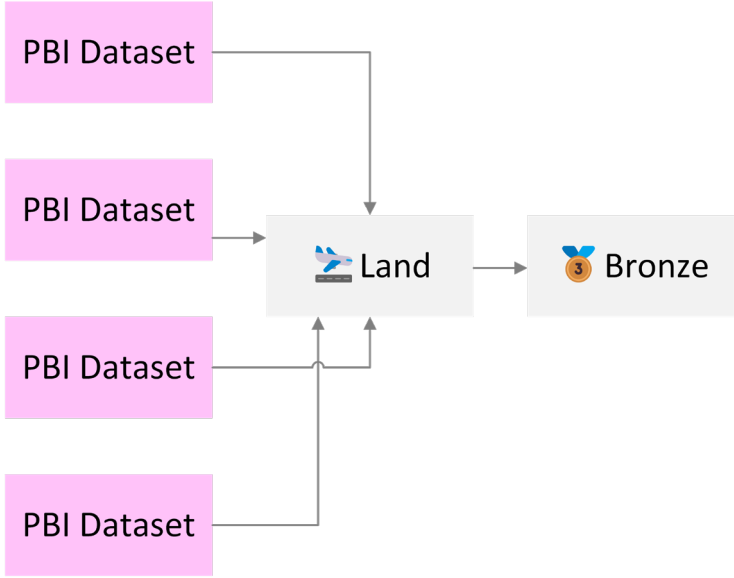
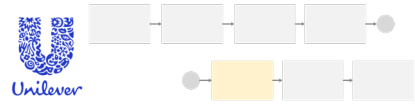
Run Flow Logic/Order

Refresh Related PBI Dashboard

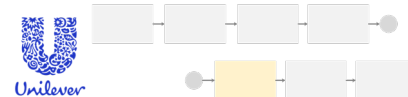


PBI PREMIUM CAPACITY

USAGE STATISTICS



D365 INTEGRATION



D365 API – SERVICE PRINCIPAL

| API / Permissions na... | Type | Description |
|-------------------------|-----------|---------------------------------------|
| ▼ Dynamics CRM (1) | | |
| user_impersonation | Delegated | Access Common Data Service as orga... |

Power Platform admin center

⊕ Add user ↻ Refresh 📄 Manage users in Dynamics 365

Home

Environments

Advisor

Analytics

Environments > Europe Data Foundations > Settings > Users

Manage users so they can access data within their environment. This list includes users with disabled and enabled statuses. [Learn more](#)

Looking for application users? Click here to go to the [app users list](#)

To validate user permissions for specific app(s), go to [app access checker](#).



D365 INTEGRATION

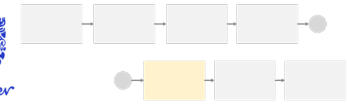
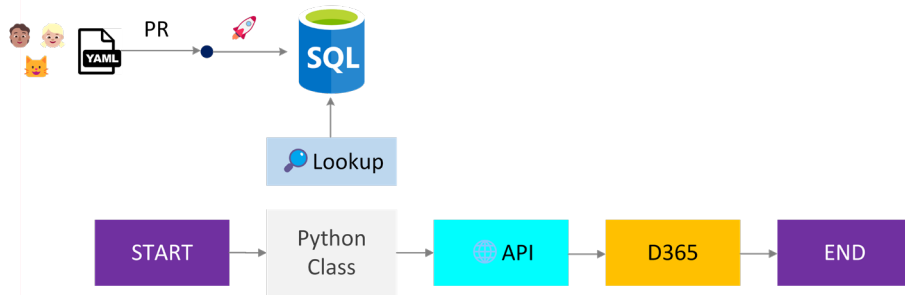


TABLE CREATION BASED ON YAML DEFINITION

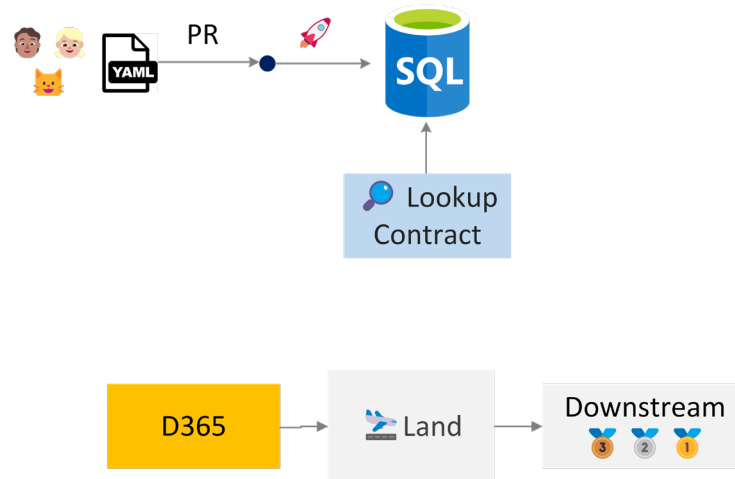
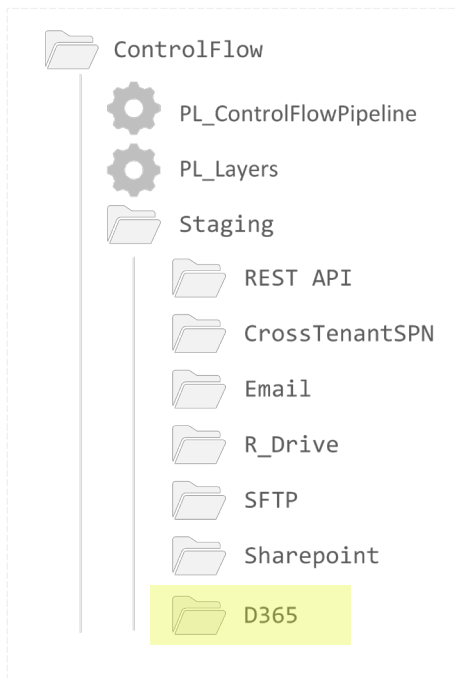
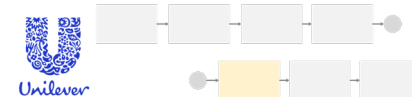


API Permissions

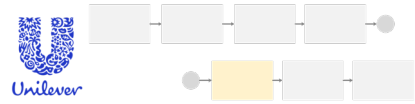
```
dv_target_yaml = ""
table:
  edf_SnOP_Target_Test:
    unique_column: ID
    column_definitions:
      - column: Timeframe
        data_type: string_100
      - column: TargetMetric
        data_type: string_100
      - column: Portfolio
        data_type: choice_portfolio
      - column: Region
        data_type: string_100
      - column: ProductCategory
        data_type: string_50 ""
```

D365 INTEGRATION

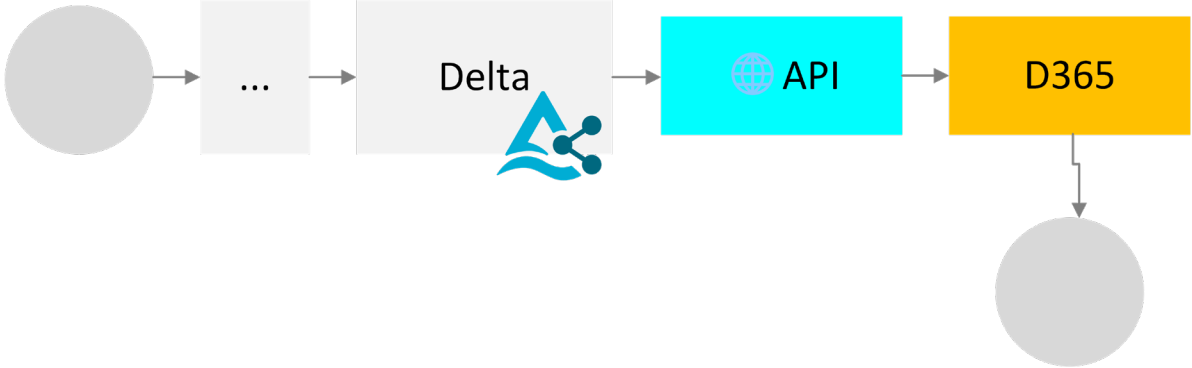
DATA FLOW – INBOUND



D365 INTEGRATION



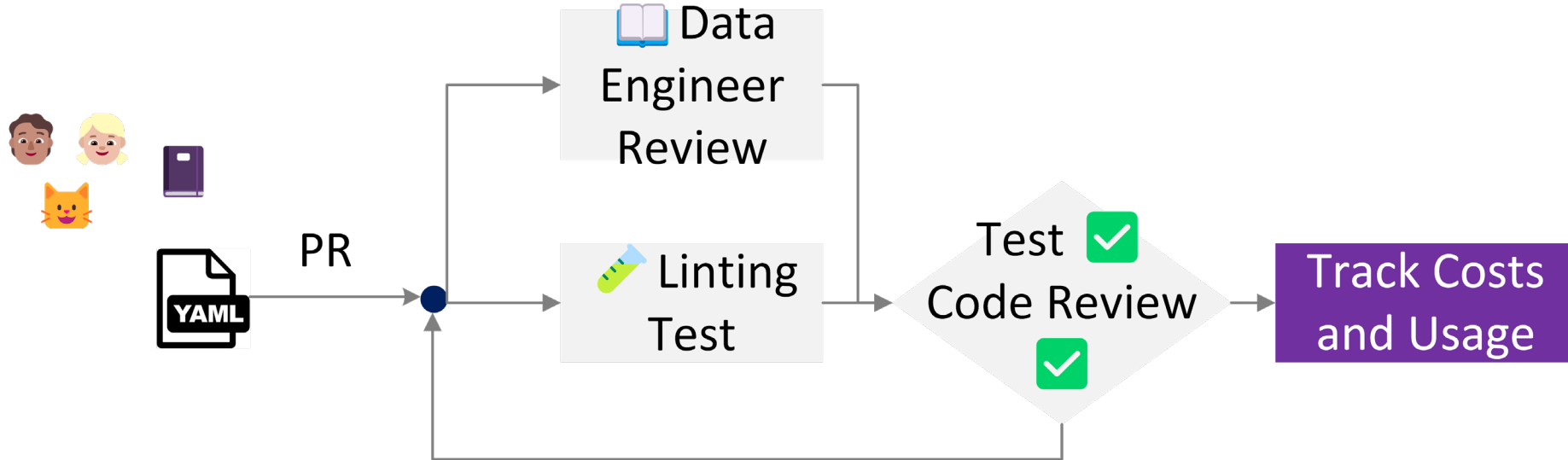
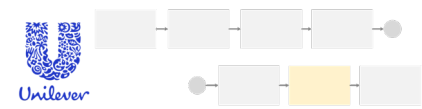
DATA FLOW – OUTBOUND



SELF SERVICE ENGINEERING

SELF SERVICE

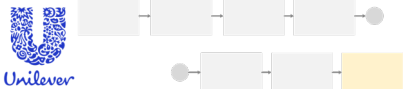
REDUCES BOTTLENECKS



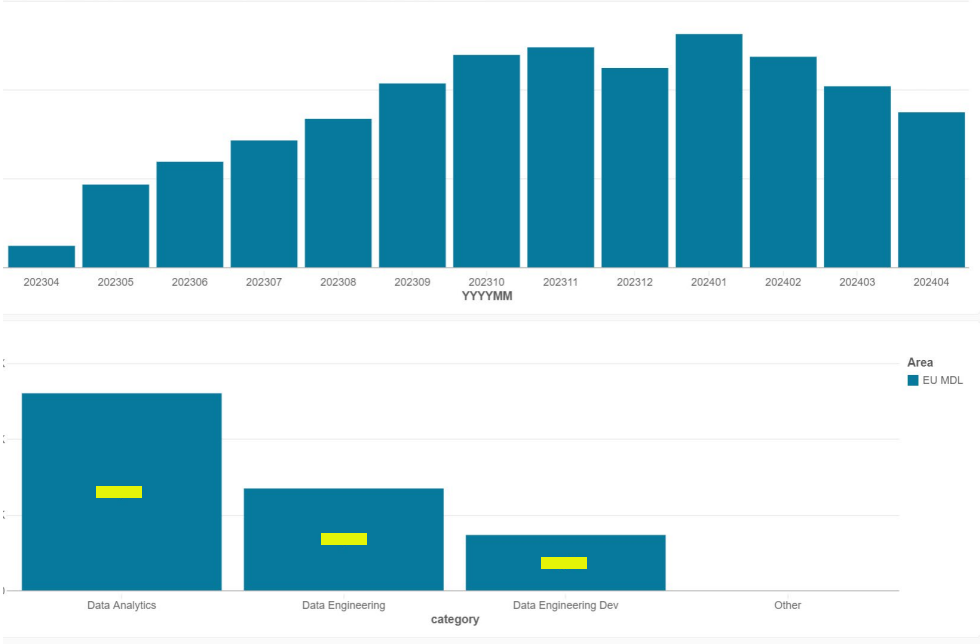
MONITORING



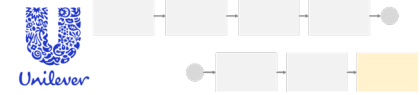
COST MONITORING



ACTIVITY BASED COSTING



USAGE MONITORING



USAGE SEGMENTED BY PERSONA AND TOOL

Role

- Business Analyst
- Data Analyst
- Data Engineer
- External
- TBD

Environment

- DEV
- PROD
- QA

Area

- ANZ PDS
- BNX MDL
- DS PDS
- EAST MDL
- EU Landing

Execution Type

- Blank User Agent
- Databricks Service / Apache Http Client
- Manual Execution

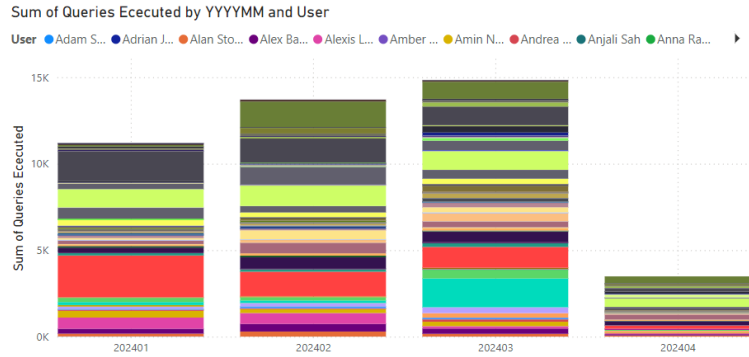
User Group

- DACH users
- TBD

User

All

☰ 🗄️ ⋮



THANK YOU!

Q&A

