# Development of a flexible, multi-stage data pipeline

## … for enhanced automation, quality control and observability

**Christian Werner, Christof Lorenz**
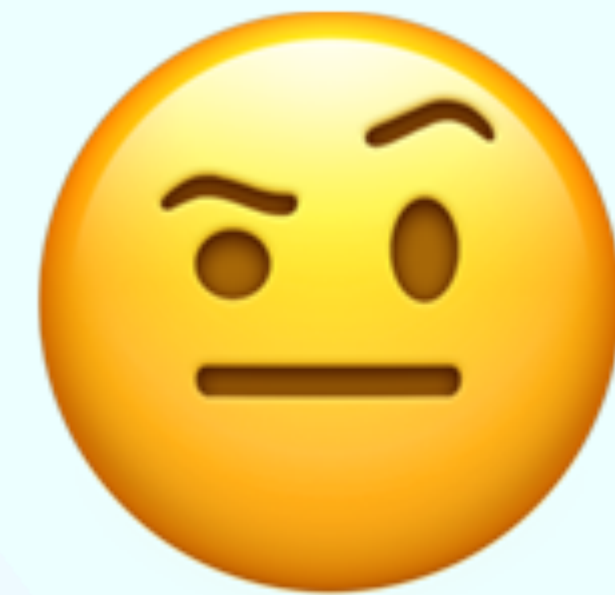**IMK-IFU, Campus Alpin, Karlsruhe Insitute of Technology**
**Contact: ✉ christian.werner@kit.edu, 🐦 @cwerner76**

**KIT**
Karlsruhe Institute of Technology

# Current Status
## … a patchwork of processes and workflows

- Heterogenuous data processing landscape

- No centralized control/ interface

- No/ limited data processing version control

- No/ limited data QA/ QC

- No centralized data catalog

# Objectives
## ... there has to be a better way

■ **Consolidate**
Centralise services, jobs and scripts into a manageable system

■ **Automate**
Create data transformation pipelines that take care of data ingestion,
auto-validation and transformation (data repositories)

■ **Increase Consistency**

A global scheduling ensures that processes run at defined times/ intervals or conditions

■ **Add Scalability**
The system should scale from few jobs to actual model processing
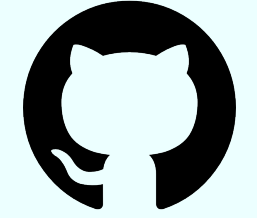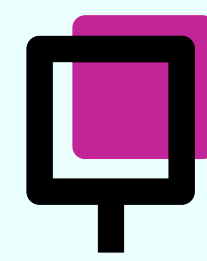(i.e. data generation for dashboards etc.)

■ **Add Observability**

System status, failure and progress should be easily observable, data owners be notified if things break or are out-of-norm
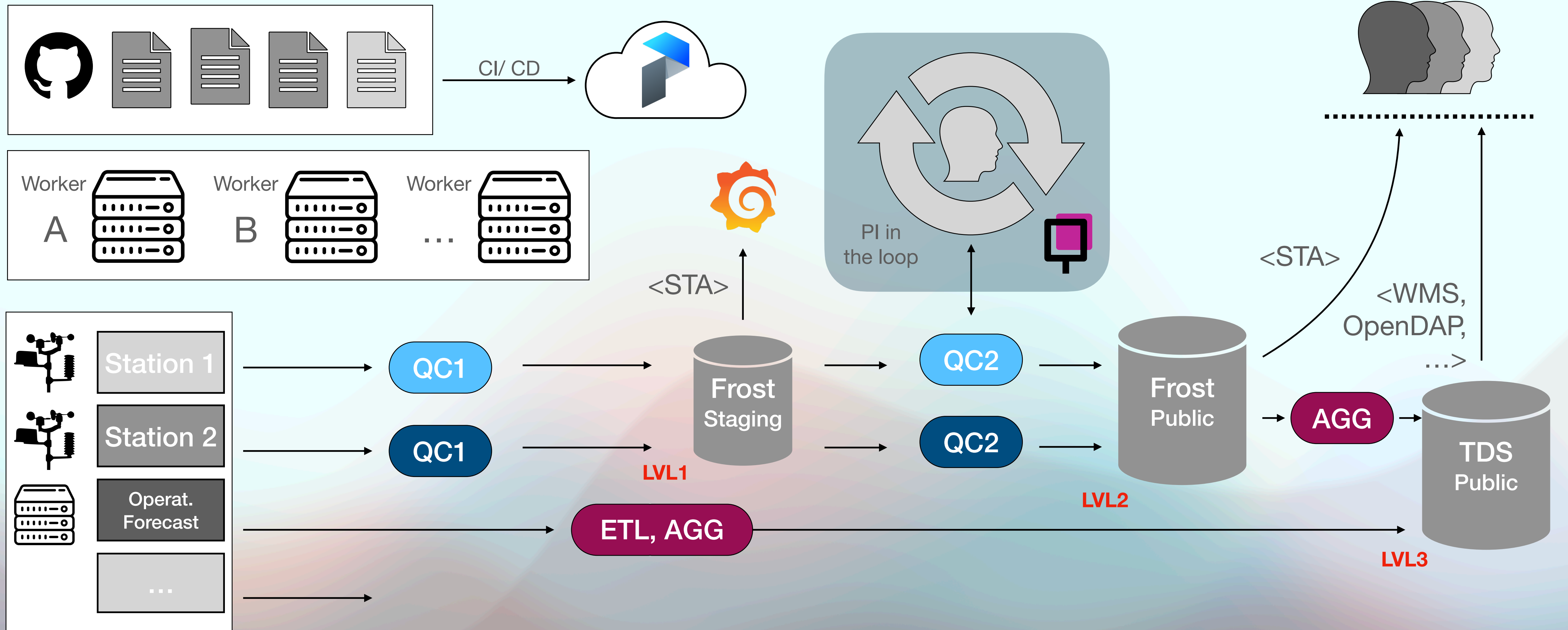
# Proposed solution
## ... develop an flexible and scalabe dataflow framework

- Open-source

- Centralized workflows with Prefect

- Workflows as code (CI/CD pipelines after change)

- Data Quality Checks with GreatExpectations and SaQC

- Data publication via STA and THREDDs

- Observability via Prefect UI, Grafana, APIs

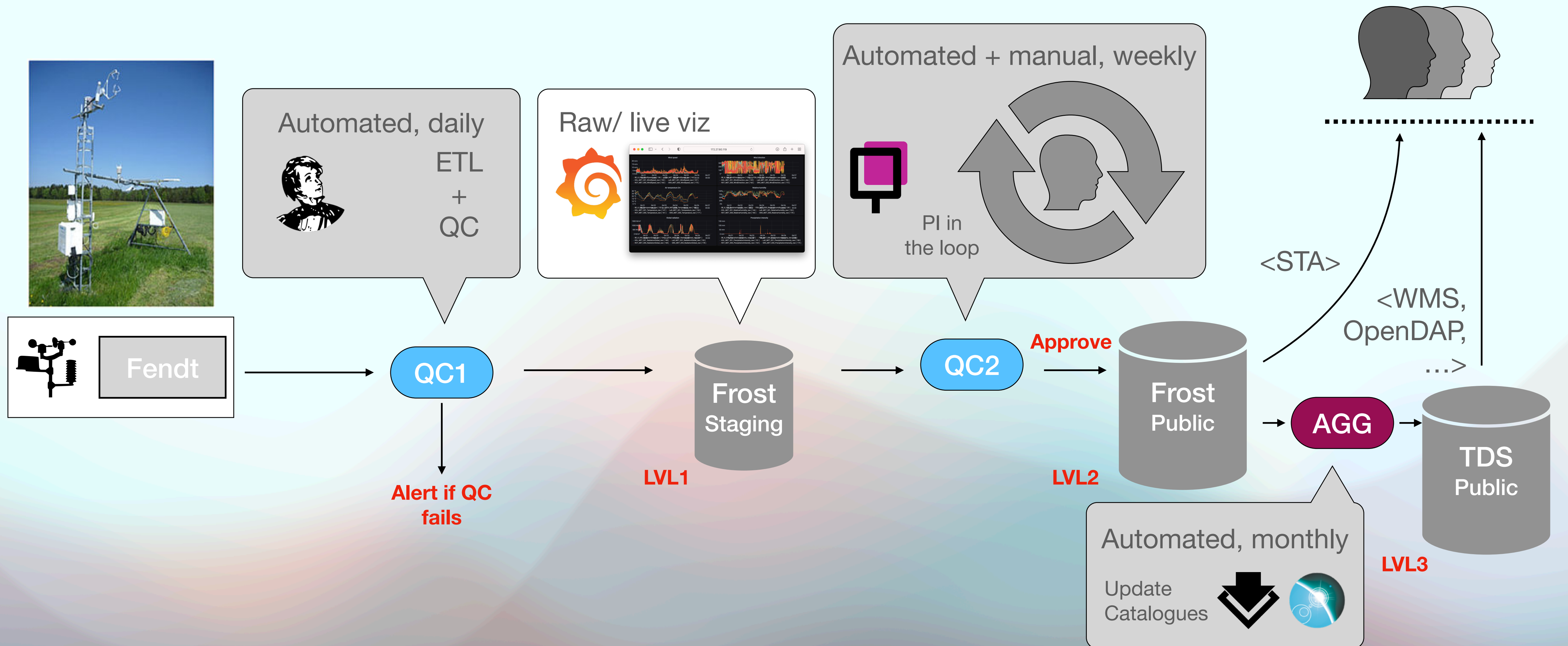# Workflow Orchestration
## … schedule and control via Prefect

# Example Workflow
## TERENO micromet ingest, QA/ QC and publication



Automated, daily
ETL + QC

Raw/ live viz

Automated + manual, weekly

PI in the loop

Fendt

QC1

Alert if QC fails

Frost Staging

LVL1

QC2

Approve

Frost Public

LVL2

AGG

Automated, monthly

Update Catalogues

LVL3

<STA>

<WMS, OpenDAP, ...>

TDS Public
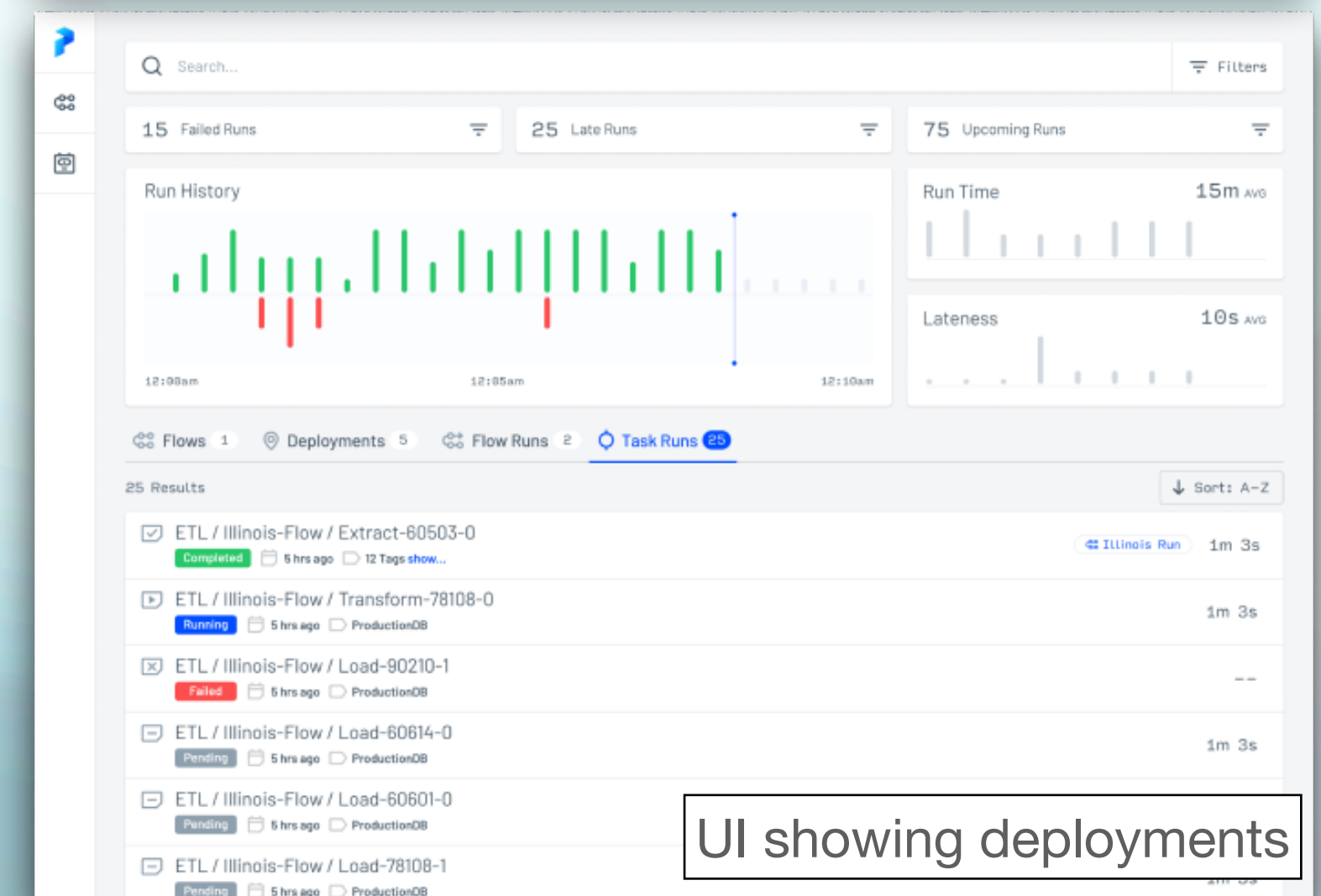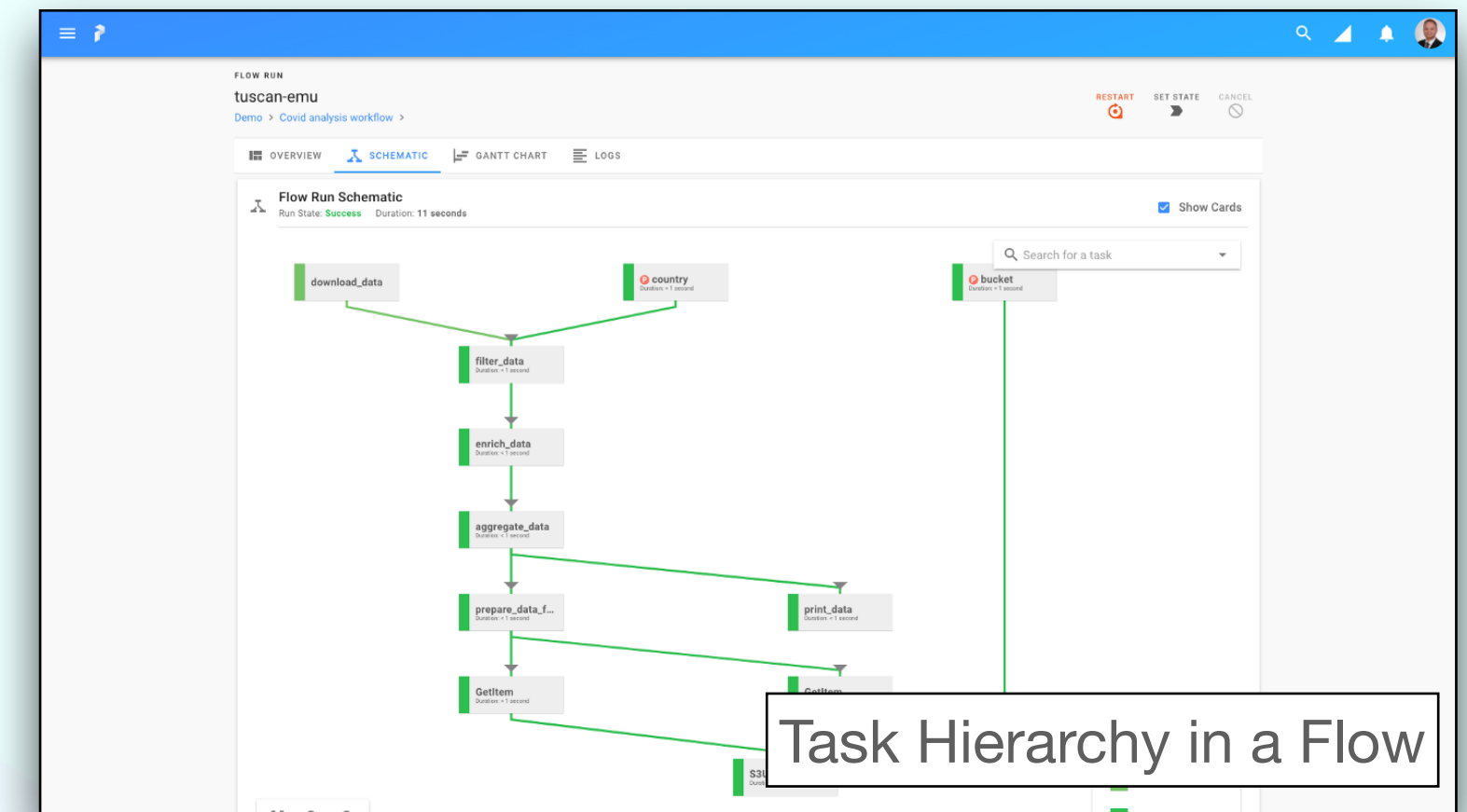
# Prefect
## Workflow definition (as code)

```python
from prefect import flow, task
from typing import List
import httpx

@task(retries=3)
def get_stars(repo: str):
    url = f"https://api.github.com/repos/{repo}"
    count = httpx.get(url).json()["stargazers_count"]
    print(f"{repo} has {count} stars!")

@flow(name="GitHub Stars")
def github_stars(repos: List[str]):
    for repo in repos:
        get_stars(repo)

# run the flow!
github_stars(["PrefectHQ/Prefect", "PrefectHQ/miter-design"])
```
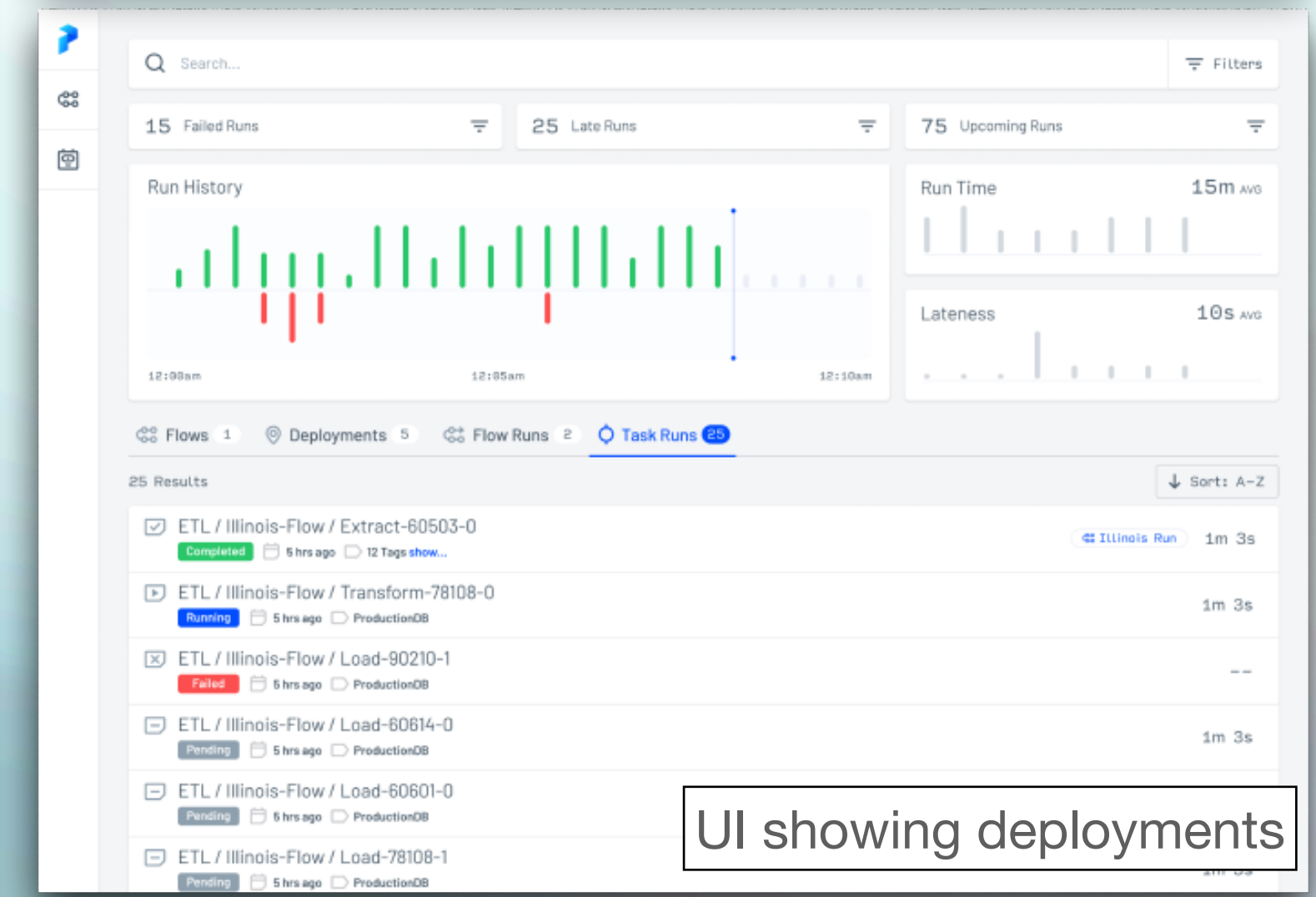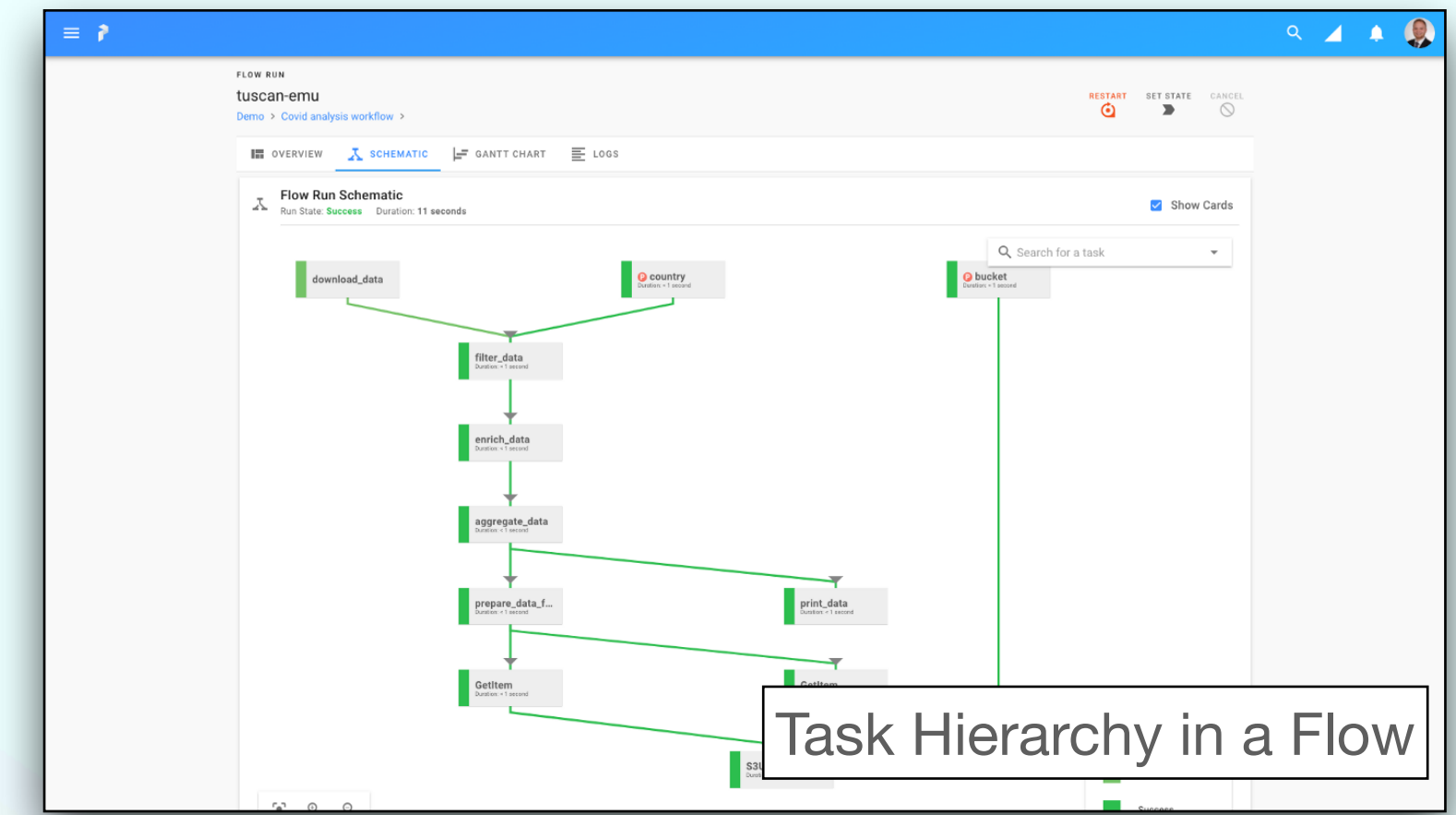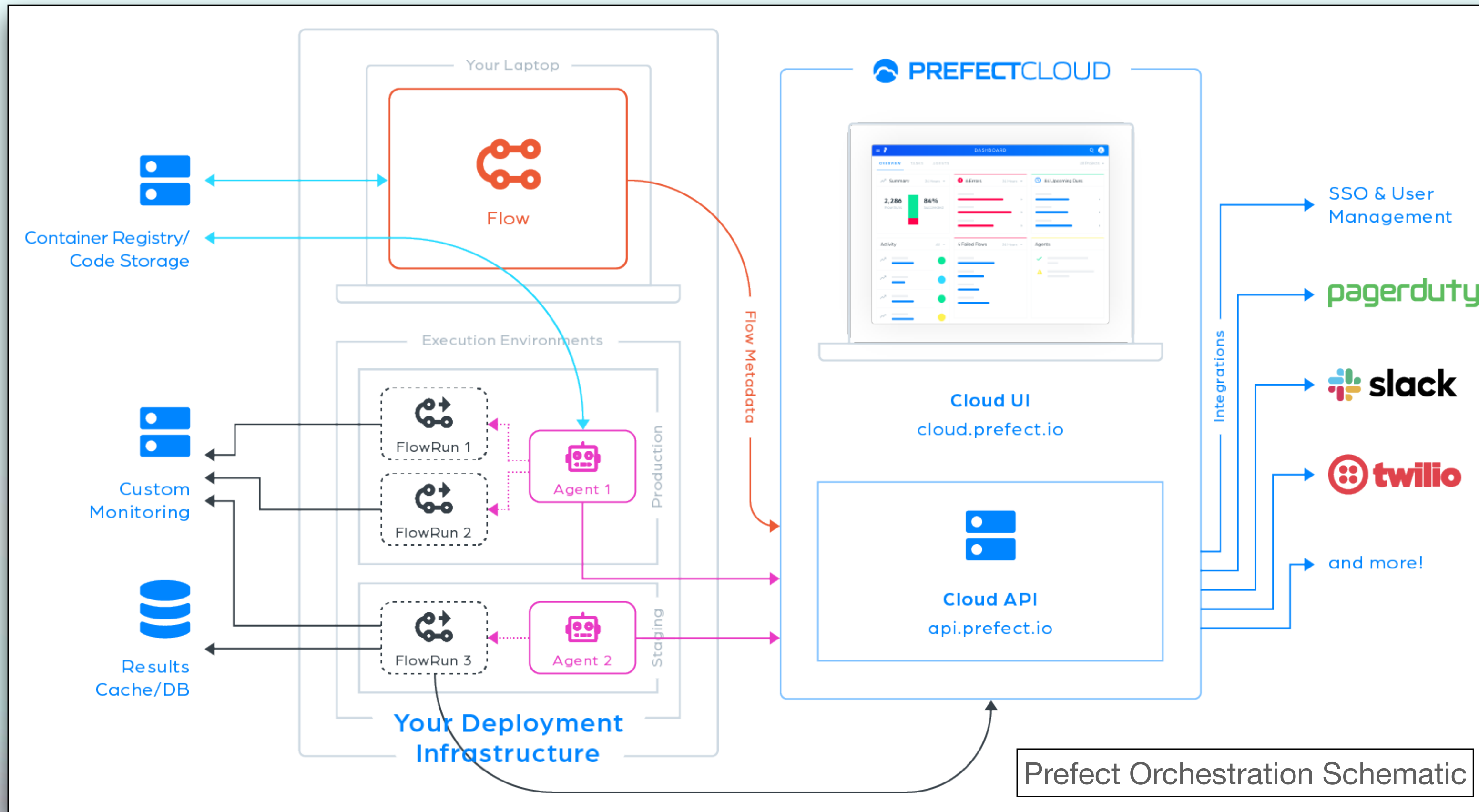
Task and flow definition as code

Task Hierarchy in a Flow

UI showing deployments

# Prefect
## Workflow orchestration



Prefect Orchestration Schematic



Task Hierarchy in a Flow
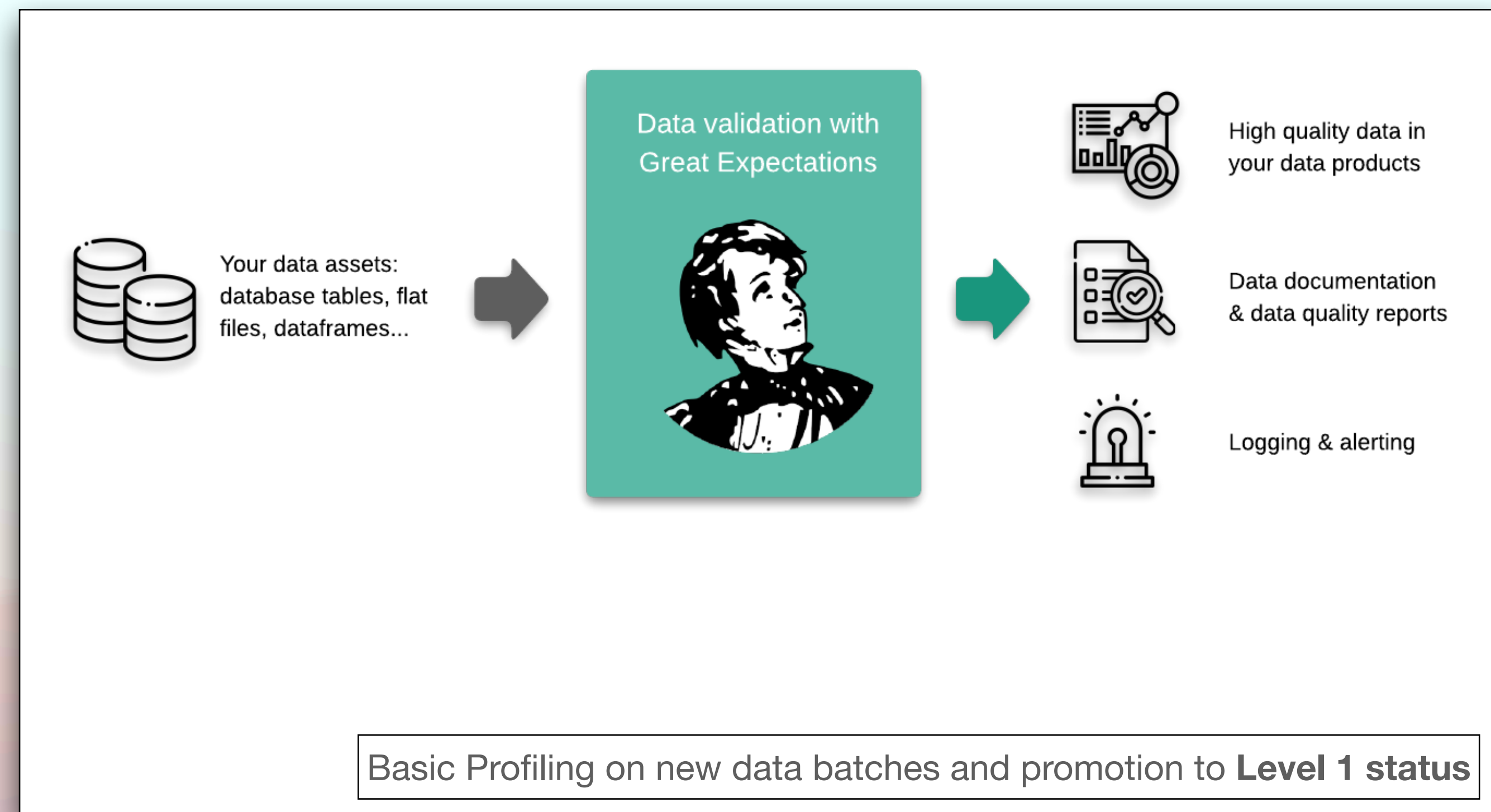


UI showing deployments

# Quality Control I
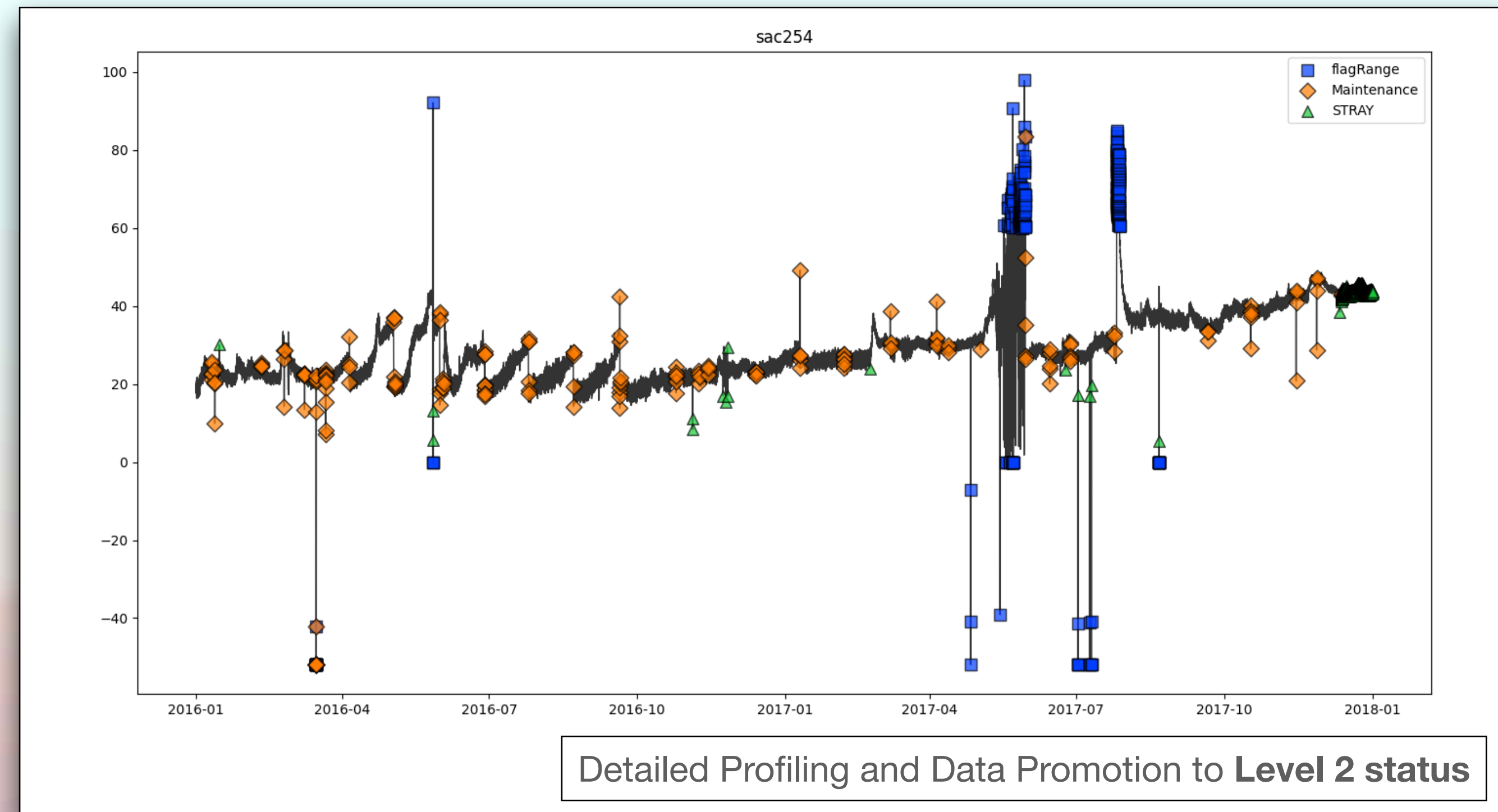## Great Expectations - Data Profiling and initial QC

- Initial automated profiling on existing valid data

- Define "expectations"

- Checks preformed at ingest time (8am), alert if missing data and other expectations are not met

- Valid data ingested into FROST staging server [LVL1]



Your data assets: database tables, flat files, dataframes...

Data validation with Great Expectations

High quality data in your data products

Data documentation & data quality reports

Logging & alerting

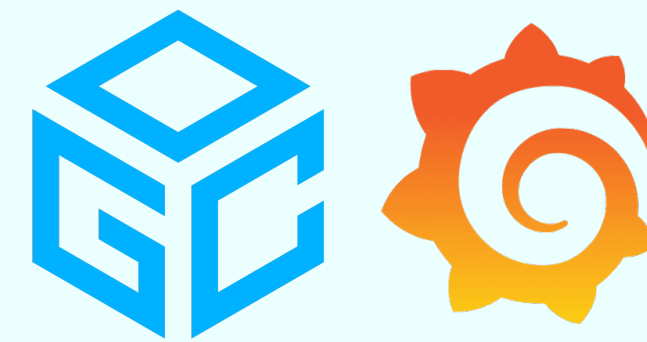Basic Profiling on new data batches and promotion to **Level 1 status**

# Quality Control II
## SaQC - measurement/ project-specific QC

- Manual intervention step at weekly interval

- Domain knowledge required

- Checks performed on time-span of data (data drift detection, outliers etc.)

- After approval data accessible via STA from 2nd FROST server [LVL2]



Detailed Profiling and Data Promotion to **Level 2 status**

# Data Access
## Sensor Things API

- Data access via STA API
  (REST or via code, i.e. stantic)

- (Grafana) Dashboards



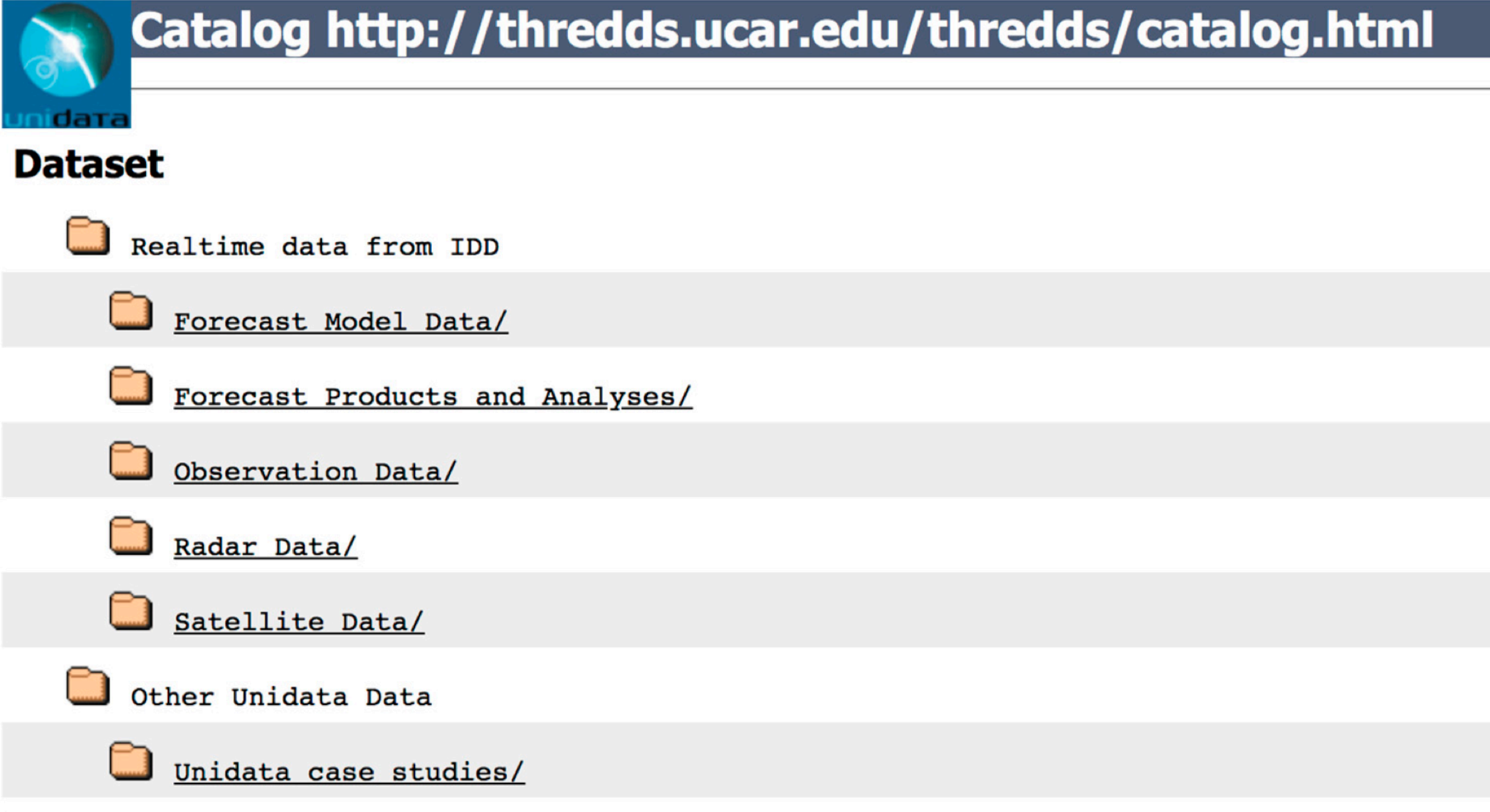Data access and live plotting of Level 1/2 data via STA

# Data Access
## THREDDS Data Server

- Aggregated datasets [LVL3] (csv, netCDF, monthly/ annual)

- Automatically expand data catalog

- Intake catalog
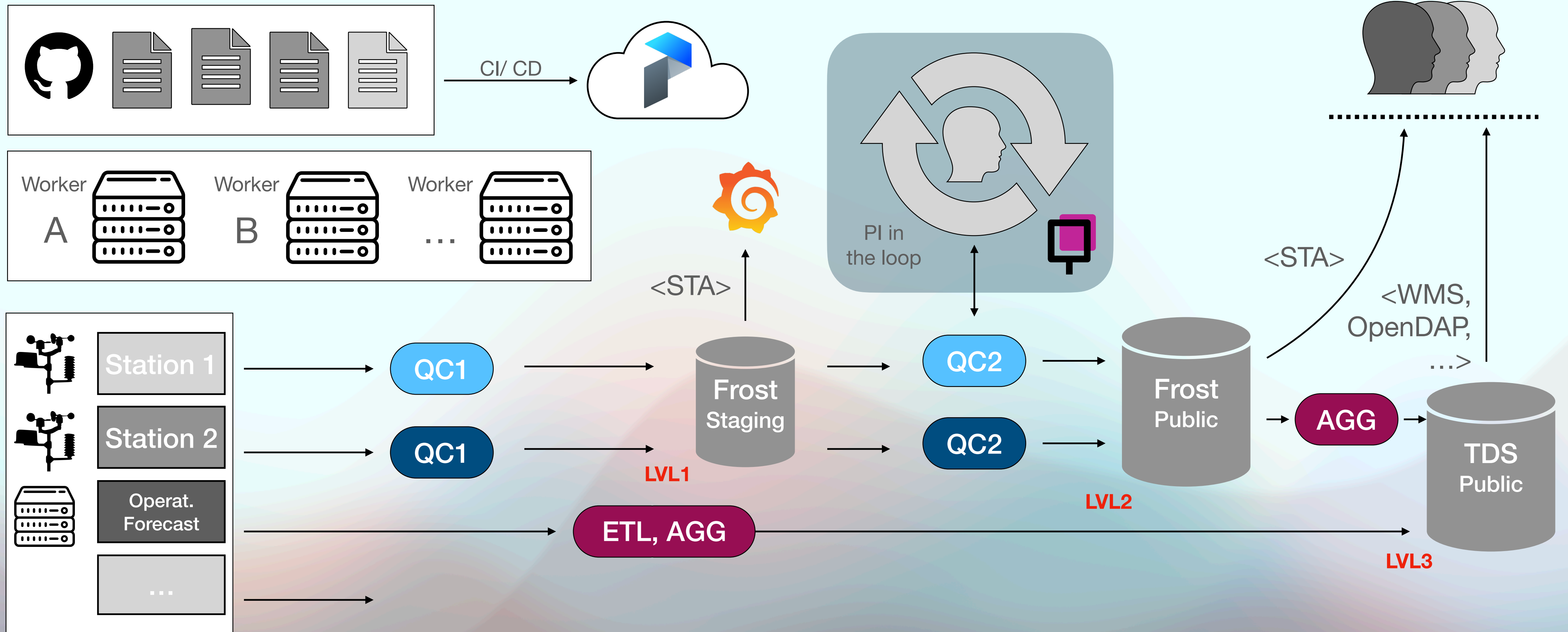
- Harvested into searchable Geoserver and/ or Data Portal

### Catalog http://thredds.ucar.edu/thredds/catalog.html

**Dataset**

- 📁 Realtime data from IDD
  - 📁 Forecast Model Data/
  - 📁 Forecast Products and Analyses/
  - 📁 Observation Data/
  - 📁 Radar Data/
  - 📁 Satellite Data/
- 📁 Other Unidata Data
  - 📁 Unidata case studies/

**THREDDS Data Server** at **Unidata** see **Info**
**THREDDS Data Server** [Version 4.6.14 - 2019-07-23T11:04:31-0600] **Documentation**

Data access through THREDDS Data Server

# Workflow Orchestration
## … schedule and control via Prefect

# References

- FROST Server: https://fraunhoferiosb.github.io/FROST-Server

- Great Expectations: https://greatexpectations.io

- Intake: https://intake.readthedocs.io

- OGC Sensor Things API: https://www.ogc.org/standards/sensorthings

- Prefect: https://prefect.io

- Stantic: https://cwerner.github.io/stantic

- System for automated Quality Control (SaQC): https://rdm-software.pages.ufz.de/saqc

- THREDDS Data Server: https://www.unidata.ucar.edu/software/tds