



D3.4 - In Situ Data Hub I

WP3 – Large Scale Demonstrators

Authors: Stavros Tekes, Yanis Nasiopoulos

Date: 12.03.20

Full Title	Promoting the international competitiveness of European Remote Sensing companies through cross-cluster collaboration			
Grant Agreement No	824478	Acronym		PARSEC
Start date	1 st May 2019	Duration		30 months
EU Project Officer	Milena Stoyanova			
Project Coordinator	Emmanuel Pajot (EARSC)			
Date of Delivery	Contractual	29.02.2020	Actual	06.03.2020
Nature	Other	Dissemination Level		Public
Lead Beneficiary	DRAXIS			
Lead Author	Stavros Tekes	Email		stavros@draxis.gr
Other authors	Giannis Nasiopoulos, Alex Economou			
Reviewer(s)	Peter Baumann (RASDAMAN)			
Keywords	API Service, EO, dockers, containers, database, architecture, hardware			

Document History

Version	Issue date	Stage	Changes	Contributor
1.0	20.02.2020	Draft	First draft	DRAXIS
1.1	27.02.2020	Draft	Merge review comments	DRAXIS
1.2	06.03.2020	Draft	Review comments	DRAXIS
1.3	12.03.2020	Final	Formatting	DRAXIS

Disclaimer

Any dissemination of results reflects only the author's view and the European Commission is not responsible for any use that may be made of the information it contains

Copyright message

© PARSEC consortium, 2020

This deliverable contains original unpublished work except where clearly indicated otherwise. Acknowledgment of previously published material and of the work of others has been made through appropriate citation, quotation or both. Reproduction is authorised provided the source is acknowledged.

Table of Contents

Table of Contents	3
List of Acronyms	3
Executive Summary	4
1 Introduction	5
2 Main Features.....	5
3 Architecture.....	7
3.1 Backend API Service	7
3.2 In Situ Data Storage Service	7
3.3 CAMS Service.....	8
4 Hardware requirements.....	9
5 Software requirements.....	10
6 Hub Expandability.....	11

List of Acronyms

API	Application Programming Interface
RESTful API	A web based HTTP Application Programming Interface
B-tree	Data structure that maintains sorted data and allows fast searches
CAMS	Copernicus Atmosphere Monitoring Service
CSV	Comma Separated Values
Datcube	Multi-dimensional Array of values
EO	Earth Observation
FTP	File Transfer Protocol
GDAL	Geospatial Data Abstraction Library (software library for handling geospatial data)
GIST	Generalized Search Tree (allows fast searches)
GIN	Generalized Inverted Index (allows fast searches)
GIS	Geographic Information System
Hash	Data structure that maintains sorted data and allows fast searches
HTTP	Hyper Text Transfer Protocol
IoT	Internet of Things
JSON	JavaScript Object Notation
Kong	Software to securely manage communication between clients and microservices
MongoDB	Cross-platform Document Oriented Database
NGINX	Software that Accelerates Content & Application Delivery
NetCDF	Network Common Data Form
PHP	General-purpose Programming Language
PostgreSQL	Advanced Open-source Database
PostGIS	Software program that adds support for geographic objects to the PostgreSQL
RDBMS	Relational Database Management System
VM	Virtual Machine (multiple PC-like machines running under one computer)
XML	Extensible Markup Language

Executive Summary

This report documents the PARSEC In Situ Data Hub which will act as a repository hosting geo-referenced field observations coming from sensor networks, independent IoT devices, citizen observatories or other online structured and unstructured data sources available.

The In Situ Data Hub will provide access to real and past time data through pre-populated measurements derived from institutional in situ datasets hosted by major stakeholders responsible for monitoring of climate, environment and different sectors of the economy (e.g. food security). Additionally the hub will collect open sensor measurements available on the internet either in structured formats (e.g. XML or JSON services, databases, csv or excel files) or as unstructured content such as sensor measurements in web page tables in a structured way.

Initially the overall platform design and architecture of the data hub is presented, followed by hardware and software requirements for the deployment of the hub.

1 Introduction

The In Situ Data Hub is a multi-source repository with automated discovery, retrieval, harmonisation and transformation services. Geospatial data from a wide range of sources is acquired and transformed into time series datacubes (multi-dimensional array of values), which are made available for use in a combined, uniform, analytical environment.

The acquired datasets are made publicly available for access and download. The retrieved information can be used for various purposes, such as

- i) the production of value-added products, particularly in combination with satellite data
- ii) validation of parameters extracted with EO techniques and algorithms
- iii) training neural networks
- iv) calibration or visualisation purposes,
- v) discovering insights based on big data analysis that can combine observations, time and geolocation information, etc
- vi) services of operational monitoring, warning and reporting to public institutions or businesses.

Currently the In Situ Data Hub provides air quality data for approximately 10 pollutants originating from:

- Official ground-based monitoring stations, with approximately 10651 daily measurements for 12000 stations and provide information as of August 2019 and onwards.
- Low cost sensors, with approximately 27394 measurements per minute, per hour and daily for 30000 sensors and provide information as of September 2019 and onwards.
- CAMS (Copernicus Atmosphere Monitoring Service) with a 3 days forecast and 10km resolution providing continuous past and present data and information on atmospheric composition.

2 Main Features

The In Situ Data Hub is an automated data acquisition and pre-processing system, which exposes API services for data provision. Geospatial data from a wide range of sources are acquired (sensor networks, independent IoT devices, citizen observatories, etc.), transformed into time series, and ingested into the database, ready for use in a combined, uniform, analytical environment.

Currently the hub integrates in-situ air quality (AQ) data, from official monitoring stations, low-cost sensors (IoT-Internet of Things networks) and air quality forecasts from the Copernicus Atmosphere Monitoring Service (CAMS).

The main objective of the hub will be to offer SMEs, a centralised, expedite access to valuable resources that is often key to realising value-added EO(Earth Observation) services.

An indicative screenshot of the In Situ Data Hub interface is presented below.

In-situ data hub

In-situ data hub is a collection of various free to use data APIs

Air Quality: Sensors

Get the collection

Get the last measurements

Get the capabilities of a sensor

Get the capabilities of a sensor

Get the historical measureme...

Air Quality: Stations

Get the stations collection

Get the last measurements

Get the capabilities of a station

Get the capabilities of a station

Get the historical measureme...

Air Quality: CAMS Forecast

Get the CAMS forecast for a ...

Air Quality: Sensors

Air Quality is a simple API allowing consumers to view various Air Quality measurements from different sources

SENSORS

Sensors is a collection of low cost sensors around the world. This API retrieves the basic info of them. We have added various query URI template parameters

GET

/airquality/sensors{?limit_records,sensor_state}

Get the collection

Example URI

GET /airquality/sensors?limit_records=&sensor_state=

URI Parameters

limit_records

number (optional) Default: unlimited

The maximum number of results to return.

sensor_state

text (optional) Default: "active"

Select active / inactive or all sensors

Request

Sensors

Show

Response

200

Show

LAST MEASUREMENTS

Gets the last measurements of the low cost sensors. We have added various query URI template parameters

GET

/airquality/sensors/last_measurements{?limit_records,sensor_state,limit_days_old}

Get the last measurements

Example URI

GET /airquality/sensors/last_measurements?limit_records=&sensor_state=&limit_days_old=

URI Parameters

limit_records

number (optional) Default: ""

The maximum number of results to return.

sensor_state

text (optional) Default: "active"

Select active / inactive or all sensors

limit_days_old

number (optional) Default: ""

Do not show measurements older than that many days

Request

Sensors

Show

Response

200

Show

Figure 1: The In Situ Data Hub user interface

3 Architecture

The In Situ Data Hub adopts an agile system architecture ensuring the seamless integration of various technological components, enabling optimal use of resources, and driving EO data exploitation and EO service provision in the future. As seen in figure 2, further below, the high-level architecture of the hub will consist of three interconnected service tiers.

- Backend API Service
- In Situ data storage Service
- CAMS Service

3.1 Backend API Service

A RESTful API (web based HTTP Application Programming Interface), based on the PHP framework Lumen (a popular general-purpose programming language), is responsible to apply air-quality intelligent algorithms on the retrieved raw pollutants. Lumen is an integrated framework that facilitates application development through built-in programming tools, allowing for a fast and ease deployment compared to other frameworks of the PHP programming language.

The design of the In Situ Data Hub source code follows the MVC development model, Model - View - Controller, which ensures the best distinction between the logic of the application and data models and the presentation of the retrieved results. It enables the development process to be separate and thus achieves improved performance in every aspect of the application's development lifecycle (Growth Speed, Debugging, Upgrades).

Additionally the NGINX Component (software that accelerates content and application delivery), improves security and acts as the web server of the API's stack. It is an open-source software for web serving, reverse proxying, caching, load balancing, media streaming, and other services.

API Gateway

The Kong gateway (software to securely manage communication between clients and microservices) is the single entry point into the In Situ Data Hub infrastructure. It handles a request by invoking multiple rules such as authentication mechanisms, throttling and aggregating the results. It can also translate between web protocols such as HTTP and WebSocket and web-unfriendly protocols that are used internally.

3.2 In Situ Data Storage Service

Postgres database is the in situ data storage, powered by the geospatial component PostGIS required for georeferenced storage of measurements coming from sensors, stations and raster data of the Copernicus Atmosphere Monitoring Services.

PostgreSQL is an open source RDBMS (Relational Database Management System) which is not developed by a single company but by a global community of users, companies, and institutions. PostgreSQL supports functions such as B-tree, hash, GiST and GiN, automations, which enable high speed search and data retrieval from the database, and additionally a wide range of predefined and user-defined data types and objects.

MongoDB

The In Situ Data Hub infrastructure generates a large number of events (i.e. logging,) which contains useful information about their operation including errors, warnings, and user's behaviour. Due to its non-sql architecture MongoDB (document oriented database) is suitable for storing log data from the multiple subsystems and provides an easy way to monitor the overall architecture.

3.3 CAMS Service

The CAMS Service acts as a standalone micro service where requests for data are submitted from the core, Backend API service, of the In Situ Data Hub. Due to the nature of data provided by CAMS (raster raw data), particular programming expertise has been adopted to achieve image processing in combination with GIS technologies deployed by the core in situ service.

The services offered by CAMS have been carefully studied, to decide the most suitable way to integrate the plethora of variables (3-days hourly forecast datasets) provided daily, and combine those datasets with the in situ data. Namely special focus has been given to the following variables PM25, O3, PM10, SO2, CO.

CAMS supports data provision through API service and FTP server (File Transfer Protocol services). As CAMS API service is not for operational use, suffering from delays and throttling (as stated in Copernicus site), the FTP server has been chosen, which offers NetCDF (or grib) raster data with the most updated forecast cycle.

Additionally historical data stored in the PostGIS database, of the in situ data storage service, has been combined to produce structured GIS queries for statistical reasons.

The CAMS service components consist of the following.

1 A standalone service which downloads NetCDF layers

A Python service which makes a connection to the FTP endpoint of CAMS and downloads the latest 3-day forecast for the parameters of our interest.

2 A post processing service which converts downloaded raw data, ready to be stored into the relational PostGIS database of the In Situ Data Hub storage.

A python service retrieves and posts processed raw data using GDAL and Numpy libraries (software libraries for handling geospatial data) in order to extract data from satellite bands and store them into the PostGIS database.

3 An API which serves Hourly forecasts for PM25, O3, PM10 , SO2 , CO

A lumen based API has been built as the interface of the micro service, capable of retrieving Lat, Lon coordinates and responding with parameters of our interest in a json readable format (open standard data and file, human-readable interchange format)

In that respect, an illustrative example of how this agile architecture helps deliver EO-based value to end-users, in key emerging industries, is presented below.

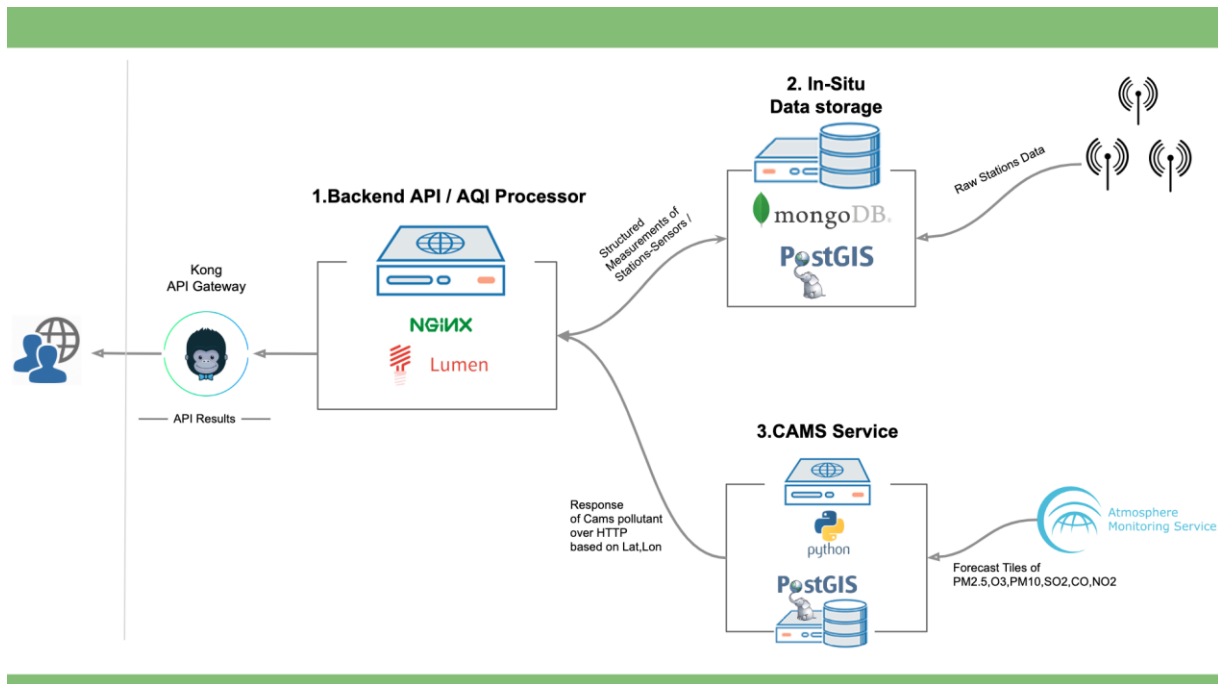


Figure 2 In-situ data hub architecture.

4 Hardware requirements

The infrastructure of the entire system is based on a single server running under Ubuntu Linux operating system and consisting of 9 Docker Containers. A Docker container image is a lightweight, standalone, executable package of software that includes everything needed to run an application. The Dockers allow the application to run anywhere, regardless of the operating system, the existence or not of a cloud infrastructure, or whether it is a physical or virtual environment. This technology enables the In Situ Data Hub to be easily migrated and deployed under any operating system and hardware infrastructure.

Through packing (“shipping”) any application into a container, acting as a separate Virtual Machine (VM), it enables the seamless integration of all the necessary dependencies to operate the service it represents, anywhere and under any operating system in exactly the same way.

The table below depicts the hardware needs of the system and the software running.

Table 1: Characteristics and Specifications of the infrastructure

Characteristics	Specifications
Operating System	Linux Ubuntu 18.04 (or higher)
Required Installed Software	Docker
CPU	4 Cores
RAM	8GB
HDD	300GB

5 Software requirements

Container based software architecture

The In Situ Data Hub follows a container based software architecture where docker images, containers, networks, volumes and other components, “ships” each subsystem making up the overall architecture acting as a unique service. Containers, working just like small-scale virtual machines (VMs), but in a far more specific and granular way, isolate each subsystem and its dependencies—all of the external software libraries the app requires to run—both from the underlying operating system and from other containers.

Some of the major advantages of Dockers and containers are listed below:

- Enables more efficient use of system resources
- Allows for faster software delivery cycles
- Enables application portability
- Reduces dependency issues

The below table along with the server-docker schema provided, depicts the different containers of the In Situ Data Hub infrastructure.

Container	Subsystems
Nginx-proxy	Proxy server filters and forwards requests to related containers
Data-api-kong	Container of API gateway (Kong)
Data-api-kong-db	Container of API gateway(Kong) Database
Data-API-environment	Codebase of Core/AQ API
Cams-API	Codebase of CAMS microservice
Mongo-logs	Container of MongoDB logging database
AQ-DB	Postgres/PostGIS Database container of In-situ data
Cams-DB	Postgres/PostGIS database container for Copernicus netcdf
Letsencrypt	Container that issues SSL certificates automatically
Docker-gen	Auto deployment container

Table 2: In Situ Data Hub container based software architecture.

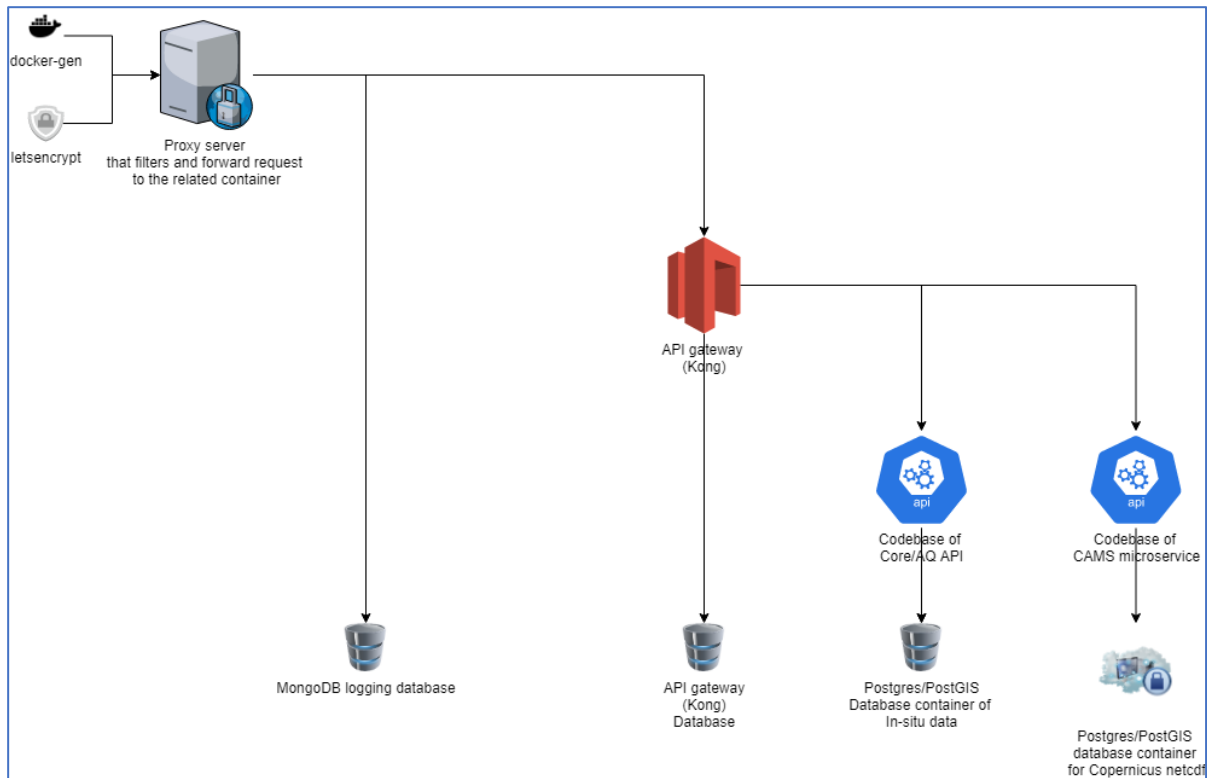


Figure 3: Server-docker schema

Scheduled data retrieval

The In Situ Data Hub contains crawlers which parse structured contents, extracting measurements with relevant information, based on predefined guide rules. Periodically scheduled cron-jobs (software used for scheduling tasks to run on the server) will trigger the crawling mechanism enabling any new data to be crawled or polled by the In Situ Data Hub. In case new data exists, it will be grabbed and ingested in the system and then forwarded to the content parser which will extract the information and store it in the database.

6 Hub Expandability

The In Situ Data Hub architecture is developed such that to allow future integration and expandability of data provision services. The software infrastructure of the hub will be able to accommodate any datasets which will arise from the discussions and needs of the beneficiaries.

Following the analysis of data needs for the PARSEC beneficiaries, a careful evaluation of the possible sources of information will be conducted. Datasets of major stakeholders, open sensor measurements available on the internet, and information from web page tables will be retrieved, parsed and ingested into the In Situ Data Hub in a structured way.

All recordsets of the hub will be expanded and enriched with new sources of information providing access to data on different domains, from air-quality, to energy, and economy.



Our Partners



This project has received funding from the European Union's Horizon 2020 research and innovation programme under grant agreement No 824478.