# Large-scale EXecution for Industry & Society

**Deliverable D7.1**

## Architectural Requirements and System Design for Interchange of Weather & Climate Model Output Between HPC and Cloud Environments

| DELIVERABLE ID \| TITLE | D7.1 \| Architectural requirements and system design for interchange of weather & climate model output between HPC and Cloud environments |
|---|---|
| RESPONSIBLE AUTHOR | James Hawkes (ECMWF) |
| WORKPACKAGE ID \| TITLE | WP7 \| Weather and Climate Large-scale Pilot |
| WORKPACKAGE LEADER | CIMA |
| DATE OF DELIVERY (CONTRACTUAL) | 31/03/2019 (M03) |
| DATE OF DELIVERY (SUBMITTED) | 01/04/2019 (M04) |
| VERSION \| STATUS | V1.1 \| Final |
| TYPE OF DELIVERABLE | R (Report) |
| DISSEMINATION LEVEL | PU (Public) |
| AUTHORS (PARTNER) | James Hawkes (ECMWF), Tiago Quintino (ECMWF) |
| INTERNAL REVIEW | Danijel Schorlemmer (GFZ), Tomáš Martinovič (IT4I) |

## DOCUMENT VERSION

| VERSION | MODIFICATION(S) | DATE | AUTHOR(S) |
|---------|-----------------|------|-----------|
| **0.1** | First Draft | 27/02/2019 | James Hawkes (ECMWF); Tiago Quintino (ECMWF) |
| **0.2** | Feedback from WP7 partners | 14/03/2019 | James Hawkes (ECMWF) |
| **1.0** | Feedback from internal review | 26/03/2019 | James Hawkes (ECMWF) |
| **1.1** | Final formatting, review | 31/03/2019 | Kateřina Slaninová (IT4I); Jan Martinovič (IT4I) |

## GLOSSARY

| | |
|---|---|
| **API** | Application Programming Interface |
| **WCDA** | Weather and Climate Data API |
| **FDB5** | Fields Database (v5) |
| **REST** | Representational State Transfer |
| **HRES** | High-Resolution (Deterministic Forecast) |
| **ENS** | Ensemble Forecast |
| **GRIB** | Gridded Binary (Data Format) |
| **NetCDF** | Network Common Data Format |
| **JSON** | JavaScript Object Notation |
| **ASCII** | American Standard Code for Information Interchange |
| **ODB2** | A data format for observations |
| **BUFR** | Binary Universal Form for the Representation of meteorological data |
| **HDF** | Hierarchical Data Format |
| **XML** | Extensible Markup Language |
| **Sensor-ML** | Sensor Model Language |
| **Water-ML** | Water Model Language |
| **AMQP** | Advanced Message Queuing Protocol |
| **AAI** | [LEXIS] Authorization & Authentication Infrastructure |
| **DDI** | [LEXIS] Distributed Data Infrastructure |

Other acronyms introduced in-place.

# TABLE OF CONTENTS

## LIST OF TABLES

## LIST OF FIGURES

ment>

# EXECUTIVE SUMMARY

The Weather and Climate Data API (WCDA) is the data management layer for curated weather & climate data in LEXIS. It is responsible for storing and organizing weather observations from a variety of sources (including in-situ unstructured observations), as well as numerical weather prediction outputs and intermediate weather & climate data. This report collects all the requirements from the LEXIS project in terms of curated weather & climate data storage and allows the preliminary design of the WCDA. It depends primarily upon the workflows defined by the Weather and Climate Large-scale Pilot (WP7).

### Position of the deliverable in the whole project context

This deliverable (D7.1) is part of Task 7.1 entitled "Develop Weather and Climate Data API", which is the first task in WP7 (Weather and Climate Large-scale Pilot). It is a precursor to D7.5 [1], which is the first release and test-bed deployment of the WCDA.

### Description of the deliverable

The deliverable begins with an introduction to the WCDA and gives more detail into the overall requirements of data curation, storage and access. Section 2 details the requirements collected from CIMA, ECMWF, ITHACA, NUM and TESEO in terms of what data they wish to provide and receive from the WCDA. Section 3 describes the preliminary design for the WCDA, including the REST interface, the preliminary system design, and the process of data curation. The main challenges in implementing the WCDA and meeting the requirements provided by this deliverable will conclude this report.

# 1 INTRODUCTION

The purpose of the Weather and Climate Data API (WCDA) is to facilitate orderly and efficient interchange of weather & climate data between the various weather and climate applications in LEXIS.



**Figure 1 Data flow into and out of the WCDA**

The WCDA will be a highly-managed, curated data storage system which is specifically designed to only handle organized weather & climate data. It will take advantage of strict data curation and indexing rules to collocate data efficiently. The WCDA is similar to the LEXIS Distributed Data Infrastructure (DDI), coming from WP3, in that it provides a means to connect data between different applications running at different geographical locations in the LEXIS project. The DDI, however, focuses on generic data storage, which is much more flexible, but also less performant for large weather datasets. The applications in the weather and climate large-scale pilot will take advantage of both the WCDA and DDI, according to their data requirements, to enable large-scale workflows (Figure 1).

The first stage of designing the WCDA is to determine what data it is required to store. For the purposes of scale and efficiency, it is important to understand the sizes of data, the frequency at which they are stored/retrieved, and the location (e.g. from ECMWF, LRZ, IT4I, etc.) where it is accessed from. For the purposes of data curation, it is important to know what the data is, including what format it is in and what scientific representation it has.

This report focuses on the collection of these requirements and the preliminary design of the WCDA based on this information. The requirements are mostly driven by weather partners' models -- but also by observational data sources that are discussed further in deliverable D7.2 [2].

## 2  ARCHITECTURAL REQUIREMENTS OF THE WCDA

This section describes the requirements of the many applications which will access the WCDA, to store and retrieve weather & climate data. These applications are described as "use-cases" of the WCDA.

### 2.1  USE CASE: PROVISION OF IFS GLOBAL FORECASTS [ECMWF]

ECMWF's main global forecasts are run twice per day, and their output will be made available to the LEXIS project via the WCDA. From there, various portions of the data can be used as initial or boundary conditions for regional models.

| APPLICATION LOCATION | |
|---|---|
| Where will you access the WCDA from? | From ECMWF's data centre |
| **APPLICATION INPUTS** | |
| What data will be **requested from** the WCDA as an input to the application? | N/A |
| At what resolution and frequency? If known, what is the size of the data? | N/A |
| In a specific data format? Which format(s)? | N/A |
| Does this application require notification of new data availability? | N/A |
| **APPLICATION OUTPUTS** | |
| What data will be **inserted in/provided** to the WCDA as an output from the application? | Any parameters from ECMWFs global HRES/ENS forecasts will be provided to the WCDA, as per ECMWFs data catalogue[1], as required. |
| At what resolution and frequency? If known, what is the size of the data? | 0.1degree lat/lon spatial resolution over 137 vertical levels with hourly temporal resolution, every cycle (12 hours). A subset of the total data output (3.4TB per day) will be provided -- estimated 100GB-1TB per day. |
| In a specific data format? Which format(s)? | GRIB |

**Table 1 Requirements - USE CASE: Provision of IFS global forecasts**

---

[1] https://www.ecmwf.int/en/forecasts/datasets/catalogue-ecmwf-real-time-products

## 2.2 USE CASE: PROVISION OF MOJI MOBILE PHONE SENSOR DATA [ECMWF]

For the duration of the LEXIS project, Moji, a Chinese weather app provider, will have a partnership with ECMWF to provide barometric data from mobile phone sensors around the world. This dataset is described further in deliverable D7.2 [2].

| APPLICATION LOCATION | |
| --- | --- |
| Where will you access the WCDA from? | Via ECMWF |
| **APPLICATION INPUTS** | |
| What data will be **requested from** the WCDA as an input to the application? | N/A |
| At what resolution and frequency? If known, what is the size of the data? | N/A |
| In a specific data format? Which format(s)? | N/A |
| Does this application require notification of new data availability? | N/A |
| **APPLICATION OUTPUTS** | |
| What data will be **inserted in/provided** to the WCDA as an output from the application? | Barometric data from mobile phones, along with latitude, longitude and time-stamp. A data point is generated every minute by every user who has the Moji weather app open. The data will be anonymized by Moji. |
| At what resolution and frequency? If known, what is the size of the data? | Approximately 10 million data points per day / 400MB per day. |
| In a specific data format? Which format(s)? | JSON, to be converted to BUFR/ODB2 for WCDA |

**Table 2 Requirements - USE CASE: Provision of Moji mobile phone sensor data**

## 2.3 USE CASE: PROVISION OF CITIZEN WEATHER STATION DATA [CIMA]

CIMA has developed a partnership with *meteonetwork*[2]*,* a citizen-scientist project to collect meteorological observations from personal weather stations across Europe. These observations, including temperature, humidity, wind and precipitation measurements, will be used by CIMA in their WRF model. This dataset is described further in deliverable D7.2 [2].

| APPLICATION LOCATION | |
|---|---|
| Where will you access the WCDA from? | Via ECMWF |
| **APPLICATION INPUTS** | |
| What data will be **requested from** the WCDA as an input to the application? | N/A |
| At what resolution and frequency? If known, what is the size of the data? | N/A |
| In a specific data format? Which format(s)? | N/A |
| Does this application require notification of new data availability? | N/A |
| **APPLICATION OUTPUTS** | |
| What data will be **inserted in/provided** to the WCDA as an output from the application? | Temperature, humidity, wind and precipitation parameters from personal weather stations |
| At what resolution and frequency? If known, what is the size of the data? | Up to 2000 stations distributed over Europe, with approximately hourly temporal resolution. |
| In a specific data format? Which format(s)? | JSON or XML, to be converted to a curated format (e.g. ODB2/BUFR). |

**Table 3 Requirements - USE CASE: Provision of citizen weather station data**

---

[2] www.meteonetwork.it[/rete/livemap]

## 2.4 USE CASE: PROVISION OF TESEO IOT DATA [TESEO]

TESEO is a system integrator who will facilitate gathering of IoT data from a variety of sources via IoT gateways. These gateways will provide meteorologically-meaningful IoT data from a range of devices. This dataset is described further in D7.2 [2].

| APPLICATION LOCATION | |
|---|---|
| Where will you access the WCDA from? | Via ECMWF |
| **APPLICATION INPUTS** | |
| What data will be **requested from** the WCDA as an input to the application? | N/A |
| At what resolution and frequency? If known, what is the size of the data? | N/A |
| In a specific data format? Which format(s)? | N/A |
| Does this application require notification of new data availability? | N/A |
| **APPLICATION OUTPUTS** | |
| What data will be **inserted in/provided** to the WCDA as an output from the application? | Temperature, humidity, wind and precipitation (and possibly others relating to air quality) parameters from IoT stations |
| At what resolution and frequency? If known, what is the size of the data? | Temporal resolution is very flexible, depending on requirements of the workflow -- every 15 minutes is sufficient. Maximum size 20MB per day. |
| In a specific data format? Which format(s)? | JSON or equivalent, to be converted to a curated format (e.g. ODB2/BUFR). |

**Table 4 Requirements - USE CASE: Provision of TESEO IoT data**

## 2.5 USE CASE: WRF MODEL [CIMA]

The WRF (Weather Research & Forecast) model is the regional weather forecast model run by CIMA, taking initial and boundary conditions from ECMWFs global weather forecast (IFS), combined with additional observations, to feed into various weather and climate pilot activities. CIMA is currently running two different WRF model instances for operational purposes: WRF-1.5km Open Loop over Italy and WRF-2.5km 3DVAR currently over north-central Italy -- but to be extended to the entire country[3]. Currently these two WRF model instances drive RISICO and Continuum model instances. Additional WRF model instances will be deployed to support additional socio-economic impact models (e.g. by NUM).

| APPLICATION LOCATION | |
|---|---|
| Where will you access the WCDA from? | CIMA/LRZ/IT4I |
| **APPLICATION INPUTS** | |
| What data will be **requested from** the WCDA as an input to the application? | (a) IFS initial and boundary condition data<br>(b) Local in situ observations (citizen scientist)<br>(c) Satellite data (ESA Sentinel data) based on CIMA and LINKs on-going work in the ESA project[4] |
| At what resolution and frequency? If known, what is the size of the data? | (a) IFS 0.1 degree resolution, approximately 1GB. Every 1-3 hours temporal resolution.<br>(b) All stations at hourly temporal resolution<br>(c) Sentinel data at 1 km grid spacing and depending on the satellite revisiting time |
| In a specific data format? Which format(s)? | GRIB, BUFR |
| Does this application require notification of new data availability? | Yes, to trigger the execution of a new WRF model execution cycle |
| **APPLICATION OUTPUTS** | |
| What data will be **inserted in/provided** to the WCDA as an output from the application? | 2m QV, Temperature, Potential Temperature<br>10m Wind Speed, Max. Wind Speed<br>Surface Pressure<br>Lightning Potential Index<br>Accumulated Total Cumulus Precipitation, Grid Scale Precipitation, Grid Scale Snow and Ice, Grid Scale Graupel<br>Downward Clear-Sky Short-Wave Flux at Ground Surface<br>Downward Short-Wave Flux at Ground Surface<br>Planetary Boundary Layer Height<br>Upward Moisture & Heat Flux at the Surface<br>Latent Heat Flux at the Surface<br>Max. Z-Wind Updraft & Downdraft<br>Max. Derived Radar Reflectivity |

---

[3] http://www.cimafoundation.org/cima-foundation/research-development/wrf.html

[4] http://www.cimafoundation.org/cima-foundation/projects/steam.html

| | Updraft Helicity & Max. Updraft Helicity |
| | Hourly Mean Z-Wind |
| | Max. Hail Diameter |
| | Pressure Level Data: Wind Speed, Temp, Rel. Humidity, Geopotential Height, Wind Speed, Dew Point Temperature, Water Vapour Mixing Ration |
| At what resolution and frequency? If known, what is the size of the data? | Between 0.01 and 0.03 degrees every 1 hours (for 48 hours) <br> As an example, over a region 1200x1200 km$^2$ wide, with 1.5 km grid spacing, a single hourly output file, providing 11 pressure levels, is around 400MB. Total approximately 20GB. |
| In a specific data format? Which format(s)? | NetCDF |

**Table 5 Requirements - USE CASE: WRF model**

## 2.6  USE CASE: RISICO FIRE RISK MODEL [CIMA]

RISICO (RISchio Incendi e COordinamento/Fire Risk and Coordination) is a mathematical model developed by CIMA Research Foundation to support operators in forest fire prevention activities. The system transforms the weather variables into information concerning fire risk. RISICO evaluates the impact of ignition on the scenario and returns a general indication of the areas where fire could be difficult to control.

| APPLICATION LOCATION | |
|---|---|
| Where will you access the WCDA from? | CIMA/LRZ/IT4I |
| **APPLICATION INPUTS** | |
| What data will be **requested from** the WCDA as an input to the application? | WRF model output, as above, at hourly temporal resolution: 10 m wind speed and direction; 2m temperature; rain; 2m humidity. |
| At what resolution and frequency? If known, what is the size of the data? | WRF 0.01-0.03 degree resolution every hour |
| In a specific data format? Which format(s)? | NetCDF |
| Does this application require notification of new data availability? | Yes, to trigger the execution of a new RISICO model execution cycle. |
| **APPLICATION OUTPUTS** | |
| What data will be **inserted in/provided** to the WCDA as an output from the application? | Percentile rate of spread <br> Mean rate of spread <br> Percentile fireline intensity <br> Mean fireline intensity <br> Fire Weather Index |
| At what resolution and frequency? If known, what is the size of the data? | Hourly temporal resolution; small total size. |
| In a specific data format? Which format(s)? | NetCDF |

**Table 6 Requirements - RISICO fire risk model**

## 2.7 USE CASE: CONTINUUM MODEL [CIMA]

Continuum is a hydrological model developed by CIMA Research Foundation to reproduce the flow of water within a basin, i.e. how much water passes into a given section of river or lake. Continuum is able to work both in the pre-event analysis and forecast phase and in the monitoring stage for the active control of hydrological events. The model has a reduced number of parameters and is also able to take advantage of all the information available via satellite.

| APPLICATION LOCATION | |
|---|---|
| Where will you access the WCDA from? | From CIMA |
| **APPLICATION INPUTS** | |
| What data will be **requested from** the WCDA as an input to the application? | WRF model output at hourly temporal resolution: weather data (10m wind speed and direction, 2m temperature, rain, 2m humidity, incoming solar radiation at surface) |
| At what resolution and frequency? If known, what is the size of the data? | WRF 0.01-0.03 degree resolution every 1 hour |
| In a specific data format? Which format(s)? | NetCDF |
| Does this application require notification of new data availability? | Yes, to trigger the execution of a new Continuum model execution cycle. |
| **APPLICATION OUTPUTS** | |
| What data will be **inserted in/provided** to the WCDA as an output from the application? | Discharge timeseries at a given catchment cross section |
| At what resolution and frequency? If known, what is the size of the data? | Hourly output; 10 minutes temporal resolution; small total size. |
| In a specific data format? Which format(s)? | WATER-ML 2.0 or ASCII (to be converted, or use generic DDI instead of WCDA) |

**Table 7 Requirements – USE CASE: RISICO fire risk model**

## 2.8   USE CASE: INDUSTRIAL SO2 PEAK PREVENTION [NUM]

NUM will use the weather forecast produced by CIMA to evaluate the improvements of forecasts in the framework of an environmental system to prevent SO2 peak around an industrial site.

| APPLICATION LOCATION | |
|---|---|
| Where will you access the WCDA from? | LRZ/IT4I |
| **APPLICATION INPUTS** | |
| What data will be **requested from** the WCDA as an input to the application? | Weather data: 10m U, V; 2m temperature; rainfall; downward short-wave flux at surface; PBL height; friction velocity at surface; surface roughness; heat sensible flux; surface pressure; 2m potential temperature; 2m humidity. |
| At what resolution and frequency? If known, what is the size of the data? | 3km spatial resolution over a 20km*20km domain in France (centered in 47°18'38.07''N / 2°03'57.25''O) from surface to at least 5km in vertical. At least 48 hours of forecast time with hourly resolution x 2 cycles ideally. Less than 1GB per day. |
| In a specific data format? Which format(s)? | GRIB (or ASCII, NetCDF, JSON) |
| Does this application require notification of new data availability? | Yes, in order to run in an automated operational mode. |
| **APPLICATION OUTPUTS** | |
| What data will be **inserted in/provided** to the WCDA as an output from the application? | Output is not curated data for the WCDA. May use DDI. |
| At what resolution and frequency? If known, what is the size of the data? | N/A |
| In a specific data format? Which format(s)? | N/A |

**Table 8 Requirements – USE CASE: Industrial SO2 peak prevention**

## 2.9 USE CASE: URBAN AIR QUALITY FORECASTING [NUM]

NUM will use the weather forecast produced by CIMA to evaluate the improvements of forecasts in the framework of an environmental system to forecast air quality over a city.

| APPLICATION LOCATION | |
|---|---|
| Where will you access the WCDA from? | LRZ/IT4I |
| **APPLICATION INPUTS** | |
| What data will be **requested from** the WCDA as an input to the application? | Weather data: 10m U, V; 2m temperature; rainfall; downward short-wave flux at surface; PBL height; friction velocity at surface; surface roughness; heat sensible flux; surface pressure; 2m potential temperature; 2m humidity. |
| At what resolution and frequency? If known, what is the size of the data? | 3 km spatial resolution over a 30km*30km domain in France (centered in 48°51'34.20''N / 2°20'20.20''E) from surface to at least 5km in vertical. 48 hours of forecast time at least with hourly resolution x 1 cycle (00Z). Less than 2GB per day. |
| In a specific data format? Which format(s)? | GRIB (or ASCII, NetCDF, JSON) |
| Does this application require notification of new data availability? | Yes, in order to run in an automated operational mode. |
| **APPLICATION OUTPUTS** | |
| What data will be **inserted in/provided** to the WCDA as an output from the application? | Output is not curated data for the WCDA. May use DDI. |
| At what resolution and frequency? If known, what is the size of the data? | N/A |
| In a specific data format? Which format(s)? | N/A |

**Table 9 Requirements – USE CASE: Urban air quality forecasting**

## 2.10 USE CASE: AGRICULTURAL DECISION-MAKING [NUM]

NUM will use weather analysis (or forecasts) produced by CIMA in order to evaluate the improvements of decision tools used by Limagrain.

| APPLICATION LOCATION | |
| --- | --- |
| Where will you access the WCDA from? | Limagrain via IT4I/LRZ |
| **APPLICATION INPUTS** | |
| What data will be **requested from** the WCDA as an input to the application? | Weather data: 2m U, V; altitude above sea (topography); 2m temperature; 2m relative humidity; rainfall; incident global solar radiation at surface; pressure at surface; downward short-wave flux at surface; upward short-wave flux at surface; downward long wave flux at surface; upward long wave flux at surface. |
| At what resolution and frequency? If known, what is the size of the data? | Spatial resolution of 3km over full European domain. Hourly data. Size not known. |
| In a specific data format? Which format(s)? | GRIB (or ASCII, NetCDF, JSON) |
| Does this application require notification of new data availability? | No |
| **APPLICATION OUTPUTS** | |
| What data will be **inserted in/provided** to the WCDA as an output from the application? | Output is not curated data for the WCDA. May use DDI. |
| At what resolution and frequency? If known, what is the size of the data? | N/A |
| In a specific data format? Which format(s)? | N/A |

**Table 10 Requirements - USE CASE: Agricultural decision-making**

## 2.11 USE CASE: SOCIO-ECONOMIC IMPACT ANALYSIS [ITHACA]

Weather and natural hazard prediction at a regional scale, with the ability to predict high-impact natural hazards (e.g. flash-flood, forest fires etc.), will be used to: (a) proactively trigger satellite/aerial (manned/unmanned) data acquisitions and (b) increase the accuracy of current emergency mapping products.

| APPLICATION LOCATION | |
|---|---|
| Where will you access the WCDA from? | IT4I/LRZ |
| **APPLICATION INPUTS** | |
| What data will be **requested from** the WCDA as an input to the application? | (a) weather data from WRF (e.g. wind speed and direction, temperature, rain, humidity, pressure, cloud cover) <br> (b) areas possibly affected by hydro-meteorological events (early-warning system) from the Continuum model |
| At what resolution and frequency? If known, what is the size of the data? | (a) Spatial resolution: variable (i.e. 1-50 km with global coverage). Temporal resolution (frequency): variable (every 1-6 hours with global coverage) <br> (b) Spatial resolution: 20-250 m. Temporal resolution (frequency): continuous monitoring |
| In a specific data format? Which format(s)? | (a) GRIB (or equivalent) <br> (b) NetCDF, JSON or equivalent |
| Does this application require notification of new data availability? | (a) Yes (or deliveries at fixed moments in time) <br> (b) Yes: a warning triggers a possible new acquisition and/or is used to validate and complement emergency mapping products |
| **APPLICATION OUTPUTS** | |
| What data will be **inserted in/provided** to the WCDA as an output from the application? | Output is not curated data for the WCDA. May use DDI. |
| At what resolution and frequency? If known, what is the size of the data? | N/A |
| In a specific data format? Which format(s)? | N/A |

**Table 11 Requirements - USE CASE: Socio-economic impact analysis**

## 2.12 OTHER REQUIREMENTS

### 2.12.1 Authorization and Authentication

It must be possible to authenticate users of the WCDA and control which data they may access. This is important for certain non-public datasets, especially in the context of third-party access from the LEXIS portal. Furthermore, it would be ideal to provide limited access (by size/number of requests) to different tiers of user. This aspect will be developed with consideration and assistance from WP4, specifically Task 4.3, which is responsible for the LEXIS Authorization & Authentication Infrastructure (AAI).

Providing access to API usage data for the LEXIS Portal (WP8) is also a requirement related to authentication and identity, but the specific implementation of this is for further study. This may be performed on a per-request basis, per data set basis or on the basis of total data transferred.

### 2.12.2 Forwarding to Other Data Services

The WCDA should serve as a mirror to any data that can be exposed to the LEXIS project from other sources. In the first instance, this includes ECMWF's Climate Data Store (CDS) and Meteorological Archival and Retrieval system (MARS). This may also expand to include Copernicus satellite observations and other services (such as Opera radar data) according to data licenses, although these sources can be accessed directly through their own APIs already.

### 2.12.3 Interface with the LEXIS Platform

The WCDA should interface with the LEXIS platform and optimize data locality based on the instruction of the orchestration system. This will probably require the WCDA to be responsible for data transportation itself or providing the necessary API to allow the LEXIS platform to move the data efficiently. Regardless, the WCDA will have to manage indexing of data across many locations.

### 2.12.4 Data Lifetime

The WCDA should have an automated means of managing the lifetime of data stored in the WCDA. The WCDA is primarily designed as a fast storage of temporary interchange data between various weather and climate models. The data required for these workflows should be automatically removed after a certain period (i.e. days/weeks), depending on the available storage hardware.

### 2.12.5 Data Slicing and Traversal

The WCDA may benefit its users by offering data slicing and traversal at the point of request. For example, if a user wishes to receive a geospatial subset or a time-series of data, the WCDA could perform this pre-processing before delivering the result to the user. This service may only apply to certain datasets or data formats, but could be a significant efficiency improvement where data locality cannot be guaranteed, and would improve user interaction.

## 2.13 SUMMARY OF REQUIREMENTS

- The data requests to the WCDA range from megabytes to terabytes. The WCDA should scale with the available hardware such that the largest requests can be handled; but should also take care to optimize for heterogeneous requests.
- The common subset of necessary data formats is GRIB, BUFR and ODB2 -- and these should be accommodated. NetCDF could also be supported with effort into data curation. It may be simpler for the relatively small NetCDF links to occur via another channel (not WCDA).
- A notification system is required to inform registered parties of new data availability.
- Access to the WCDA should be available in multiple locations (at least ECMWF, CIMA, IT4I, LRZ). As expected, this will necessitate a distributed data storage system which can operate on different hardware and perform or facilitate data transport based on instruction from the LEXIS orchestration platform.

## 3   PRELIMINARY DESIGN OF THE WCDA

In this section, a preliminary design for a front-end REST API and the back-end data storage system is described. The data curation process is also described. The design will be iterated several times during the design and development phase.

### 3.1   REST API

The WCDA will likely be implemented as a distributed REST API and follow industry norms in terms of GET/PUT/POST operations. Due to the sizes of data, their distributed nature, and the possibility of pre-processing by the WCDA, requests may take a considerable time to fulfil. Thus, the GET requests may be asynchronous, returning a poll-able URL where the client checks progress and finally retrieves their data.

Different collections of data (e.g. WRF output, Moji phone data, etc.) may be represented as separate end-points to the WCDA (subject to data curation). HTTP requests to each end-point will include authorization metadata in the header and meteorological metadata in the message content. Rather than an unstructured key-value pairing common to many generic object stores, data will be indexed using multiple scientifically meaningful keys (e.g. date, time, forecast type, height level and parameter ID). Endpoints will be provided to query collection contents and their data schemas. The REST API will also be responsible for the accounting and authorization requirements identified above.

In the final design, attention will be paid to existing standards such as OpenAPI[5] and Open Geospatial Consortium's Web Processing Service[6], as appropriate. A preliminary design is as follows in Table 12:

| URL | Method | Content | Result[7] |
|---|---|---|---|
| …/v1/collections | GET | | 200: A list of available collections. |
| …/v1/requests | GET | | 200: A list of outstanding requests and their status. |
| …/v1/{collection} | GET | | 200: A list or information regarding the data contents of a collection and/or the indexing schema. Filtering and paging provided as appropriate. |

---

[5] https://www.openapis.org/

[6] https://www.opengeospatial.org/standards/wps

[7] *4xx and 5xx errors, and other return codes, will also be returned as appropriate*

| …/v1/{collection}/data | GET | Index keys<br><br>```<br>{<br> "key1":"value/range",<br> "key2":"value/range",<br>  …<br>}<br>``` | 200 OK: Immediately returns the data.<br>202 Accepted: Returns a URL where the user can poll for their requested data.<br>303 See Other: Returns a URL where the user can download their requested data (useful for large data requests or forwarding). |
|---|---|---|---|
| …/v1/{collection}/data | POST/PUT | Data which can be decoded to determine indexing, and/or decorated with indexing keys. | 200 OK: PUT resource updated.<br>201 Created: PUT resource created.<br>202 Accepted: POST resource created.<br>Error if POSTing and resource already exists. |
| …/v1/{collection}/notification | POST/PUT | Index keys | 200 Accepted: Returns a URL where the user can poll for their notification. |

**Table 12 Preliminary design of REST API**

A notification system using a purely RESTful API may not be adequate, as it still places a requirement on the user to poll the WCDA. A publish-subscribe method using a message-oriented protocol such as AMQP may be more suitable.

## 3.2 SYSTEM DESIGN

ECMWF has a flexible object storage solution which is well optimized for weather & climate data, called FDB5 (Fields Database v5) [1]. Usage of FDB5 will allow the WCDA to meet the requirements of data sizes (MBs to TBs) and will perform well even with heterogeneous data sizes. It can be adapted to handle all required data formats (GRIB supported already; BUFR, ODB2 and others can be added) and has the fundamental aspects of access control.
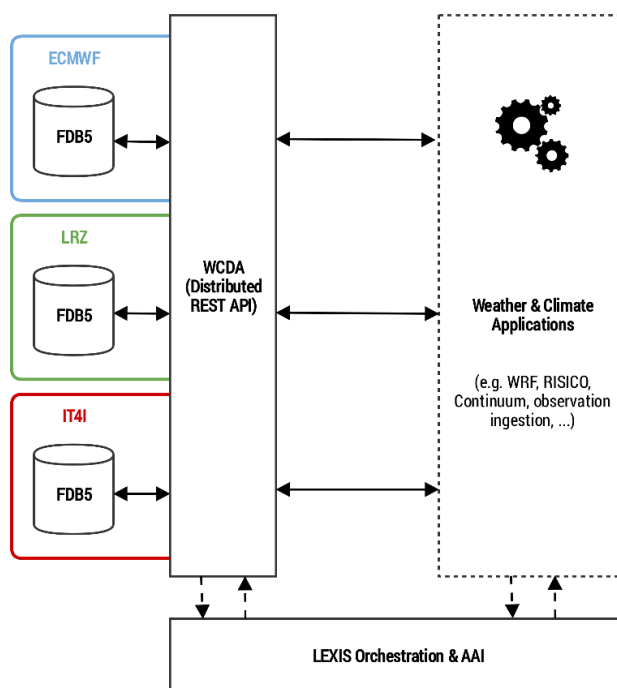
Most commercial object stores use simple key-value pairing, where the key (typically a hash) has no meaningful scientific or physical information. FDB5 uses multiple human-readable keys corresponding to scientifically-meaningful parameters to index data and this allows it to perform exceptionally well on large structured datasets typically arising from weather & climate workflows. Some data formats (such as GRIB) are self-describing, and FDB5 can use the data itself to infer the metadata keys.

The main drawback of this mechanism is that, although you can store any type of data in FDB5, the data must be curated. This is the main hurdle for most of the data that will be managed by the WCDA for LEXIS, and is discussed further in the following section. Small weather datasets noted in Section 2 as 'not curated' (typically graphs, generated reports, or other final results) can use more generic storage systems such as the LEXIS DDI system.

As of writing, the intention is to host FDB5 on each storage platform available in the LEXIS project (at ECMWF, LRZ and IT4I). FDB5 will form the back-end for the WCDA, and the aforementioned REST API will provide the front-end interface to applications (at LRZ, IT4I, CIMA, and elsewhere). There are several important developments which must be made to adapt the FDB5 for usage in the WCDA. FDB5 must have the ability to:

- Handle observation data formats (BUFR & ODB2),
- Transfer data between FDB5 peers, as per LEXIS orchestration instructions,
- Perform slicing and traversal, where possible,
- Interact with the REST API.

A preliminary system design is shown in Figure 2.



**Figure 2 Preliminary system design of the WCDA with distributed REST API and FDB5 storage. Details on LEXIS orchestration, DDI, AAI and computing are intentionally omitted as they do not impact the core WCDA design.**

## 3.3 DATA CURATION

Data curation will be one of the most significant challenges in implementing the WCDA, but will lead to significant advantages in terms of data management and performance. Data curation means that:

- The data must be described by metadata according to an agreed convention,
- The schema of metadata keys used to index the data in the WCDA must be agreed.

ECMWF has much experience with data curation for GRIB data. All the data stored in ECMWF's meteorological archive (MARS), currently over 200PB, is curated GRIB data which can be addressed using specific keys [2]. For example, a single field in the archive is addressed as:

```
class=ei, stream=oper, expver=0001, levtype=sfc, param=165.128, step=0,
time=00,
date=2019-01-01, type=an
```

The perpetuity of tools, workflows and archived data that are built upon these defined schemas means that the process of agreeing on metadata and indexing is difficult and crucial. Each of the new data types identified by this report (including ODB2/BUFR formats) will require an extensive data curation process to identify (1) the smallest non-divisible data chunk to be expressed and (2) the required metadata keys to address this data.

Some partners have identified a requirement to store NetCDF data in the WCDA. NetCDF data will require very careful data curation if it is to be included, due to the flexibility of the data format. If possible, it may be better to convert this data to GRIB or to transport this data the LEXIS DDI.

# 4 CONCLUSION

This document has focused on gathering requirements of the WCDA, driven by partners' models, observational data sources, and other needs of the consortium. A preliminary design for the front-end REST API and back-end object storage system has been presented. The initial design meets all of the primary requirements identified in terms of data management and has the scope to meet all other requirements (interaction with orchestration; authorization and authentication; data lifetime management) via the REST API.

The main challenges in implementing this design will be the development of the REST API; adaptations to the FDB5 to meet the needs of the WCDA; coupling the WCDA to the LEXIS orchestration/AAI systems; and the extensive curation of weather & climate data (especially in new formats, such as BUFR, ODB2 or NetCDF).

Work will now continue on Task 7.1, to develop the WCDA according to the requirements and preliminary design described here.

# REFERENCES

[1] LEXIS Deliverable, *D7.5 First release and test-bed deployment of Weather and Climate Data API for both in-situ unstructured observations and model data output.*

[2] LEXIS Deliverable, *D7.2 Architectural requirements and system design for interchange of in-situ unstructured weather & environmental obserations.*

[3] B. Raoult, S. Smart and T. Quintino, *A Scalable Object Store for Meteorological and Climate Data,* 2017.

[4] B. Raoult, *The Architecture of the New MARS Server,* Sixth Workshop on Meteorological Operational Systems, 1997.