

# **Open Forecast – HPC Model**

© The Open Forecast Project www.open-forecast.eu @OpenForecastEU



Co-financed by the Connecting Europe Facility of the European Union



Project	Open Forecast
Action number	2017-DE-IA-0170
Report title	Open Forecast – HPC Model
Report number	OF_M04_2020
Date of publication	15.05.2020
Revision	Rev_01
Lead Partner	HLRS
Authors	A. Schamakina, T. Schwitalla, S. Bingert, T. Bönisch
Dissemination level	<pu></pu>



## **Executive Summary**

This document describes the domain-specific HPC models and HPC workflows that will be implemented in the project.



## 1 Table of Contents

Executive Summary	2
List of Figures	4
List of Tables	5
2 Introduction	6
2.1 Project Overview	6
2.2 Project Goals	6
2.3 Purpose of this Report	7
3 High Performance Computing	7
3.1 Definition	7
3.2 Resources at HLRS	7
3.3 Difference to classical computing	8
4 Use Case I: PMFS – Particulate Matter Forecast Service 1	0
4.1 PMFS Workflow 1	5
4.1.1 Requirement analysis 1	5
4.1.2 Implementation	6
4.1.3 PMFS Case Study 1	9
5 Use Case II: AgriCOpen – Satellite Data Service for Agriculture 2	20
5.1 Requirement analysis	20
5.2 Implementation	20
6 Acknowledgements	21
7 References	!1
Appendix A: Results of the first case study 2	23



## List of Figures

Figure 1: An architecture of the Cray XC40 system
Figure 2: Domain setup of the PMFS over the Stuttgart metropolitan area. Upper row: 1250 m
resolution (left) and 250 m resolution (right). Lower panel: 50 m resolution 11
Figure 3: LUBW land cover data set prepared for the WRF simulations
Figure 4: Speedup of the WRF-CHEM simulation w.r.t. to 480 cores (left) and memory
consumption/core (right)
Figure 5: UML Activity Diagram 17
Figure 6: UML Deployment Diagram for the PMFS use case
Figure 7: Workflow to set up and operate the PMFS within Open Forecast
Figure 8: PM10 concentrations (ug/m <sup>3</sup> ) and 10-m wind velocities at 06 UTC on January 21, 2019.
Figure 9: Same as Fig. 8 but for NOx . The NO <sub>x</sub> concentrations show high values of more the 100 $\mu$ g/m <sup>3</sup>
Figure 10: Same as Fig. 8 but for SO <sub>2</sub> . The SO <sub>2</sub> concentrations appear to be accumulated in the
Neckar valley
Figure 11: Same as Figure 8 but for NH <sub>3</sub> . The values are relatively low as they are mainly related
to agriculture which is in a winter break during January 24
Figure 12: Zoom into Stuttgart downtown showing 2-m temperature and 10-m wind field. Due to
the high resolution of 50 m, the lower temperature over the elevations can be clearly identified (blueish colors)
Figure 13: Ground heat flux. Positive values indicate heat transfer from the subsurface to the
surface. Due to the applied high resolution of 50 m, the different behavior of vegetated, urban and industrial areas in terms of heat transfer is clearly visible
Figure 14: Vertical cross section of PM10 at 07:30 UTC (8:30 am local time). The center latitude
is Stuttgart main station. The Black regions denote the model terrain. It is clearly seen, that the
PM10 concentration accumulates up to 200 m above ground while it reaches almost zero above
1000 m above sea level



## **List of Tables**

able 1: Basic queues on Hazel Hen
-----------------------------------



## 2 Introduction

#### 2.1 Project Overview

The overall goal of the proposed action OPEN FORECAST is to deliver a novel Generic Service to complement the Public Open Data Digital Service Infrastructure. This HPC Open Data Forecast Service combines highly valuable, public and open data sets with supercomputing resources to produce novel relevant data products for European citizens, public authorities, economic actors, and decision makers. Analysing data sets from two innovative application areas, pollution data and agricultural data, OPEN FORECAST provides forecast services for smart farming and smart cities. Supercomputing resources are used to compute domain specific methods on large data sets to generate forecasts on urban pollution and for precision farming. The resulting data products are public and open and will be made available through the European Data Portal. Additional publication channels include domain-related data portals and APIs for the integration into stakeholder services. Furthermore, data visualization offered by OPEN FORECAST allows researchers and experts to conduct visual analyses of the data products and serve as visual communication assets towards citizens and decision makers. The whole service pipeline is "designed to be extended and to be re-used". This approach enables other European stakeholders and use-case owners from other application domains to adapt their business models to the HPC Open Data Forecast Service pipeline. OPEN FORECAST exploits, wherever possible, results from European activities. This includes the usage of the eID CEF<sup>1</sup> Building Block and the publication of project results through INSPIRE<sup>2</sup> services. The data products are enriched with metadata, compliant with standards used for the European Data Portal, and APIs are provided for seamless harvesting.

#### 2.2 Project Goals

The overall goal of the action is to combine public and open data sets with supercomputing resources to produce novel data products for European citizens, public authorities, economic actors, and decision makers. Specifically, the Action will deliver a service (the "Open Forecast service") that will take data from different sources, process it using supercomputing resources to generate open and public data products compliant with the European Data Portal.

The function of the service, which will be designed as generic as possible to potentially serve a multitude of use cases, will be exemplified through the implementation of two use cases: Particulate Matter Forecast Service (PMFS) and Satellite Data Service for Agriculture (AgriCOpen). PMFS's goal is to provide a detailed forecast for the concentration of particulate matter in the Stuttgart metropolitan area. The use case will use open sensor data from the luftdaten.info<sup>3</sup> network and public data from LGL, which will be processed with supercomputing resources. The resulting open data products will be particulate matter forecasts. AgriCOpen's aim is to provide services and products for agricultural smart farming practices based on open satellite

<sup>1</sup> https://ec.europa.eu/cefdigital/wiki/display/CEFDIGITAL/eID

<sup>&</sup>lt;sup>2</sup> https://inspire.ec.europa.eu

<sup>&</sup>lt;sup>3</sup> <u>https://luftdaten.info</u>, while writing this document it changed to https://sensor.community



imagery data like Sentinel-2<sup>4</sup>. It will use public open spatial data from the Sentinel-2 satellite mission. Data will be processed via supercomputing resources.

The resulting open data products will be agricultural forecast data (crop parameters, yields, productivity stability measures) which will be integrated into smart farming applications in the field (e.g. fertilisation or precision farming) and will be compared to ground truth data like nitrogen sensor data, yield maps, or soil maps to evaluate their applicability.

### 2.3 Purpose of this Report

The document is to define the domain-specific models and workflows. Based on those definitions the computational requirements can be derived. The necessary code and programs will be identified.

## 3 High Performance Computing

High-Performance Computing (HPC) is computing performed on computer systems with specifications that far exceed conventional computers. Such computer systems have a high-level performance. The performance of a supercomputer is commonly measured in floating-point operations per second (FLOPS). HPC is the use of distributed computing facilities for solving problems that need large computing power. Historically, supercomputers and clusters are specifically designed to support HPC applications that are developed to solve "Grand Challenge" problems in science and engineering.

#### 3.1 Definition

HPC systems are typically an aggregation of a bunch of computers, each one of which can look pretty similar to a personal computer. Each individual computer in a cluster is a node. A node has own processors, memory, disks, and an operating system. Nowadays processors typically have from two to four cores. However, cluster CPU processors can have up to 96 cores and graphical processors have hundreds of cores. The point of having a high-performance computer is so that the individual nodes can work together to solve a problem larger than any single computer can solve. The nodes need to be able to communicate in order to work together.

#### 3.2 Resources at HLRS

HLRS offers a variety of supercomputing systems reflecting the different needs of its scientific and industrial customers. At the project start, the supercomputer Hazel Hen was at the heart of the high-performance computing system infrastructure at HLRS. With a peak performance of 7.42 Petaflops, Hazel Hen was one of the most powerful HPC systems in the world (position 27 in the TOP500, July 2018). Hazel Hen entered operation in October 2015, was based on the Intel® Haswell Processor and the Cray Aries network technologies and was designed for sustained application performance and high scalability. Hazel Hen had 7712 compute nodes with 185,088 compute cores, 128 GB of memory per node and the disk capacity is around 10 PB.

<sup>&</sup>lt;sup>4</sup> https://www.esa.int/Our\_Activities/Observing\_the\_Earth/Copernicus/Sentinel-2



In February 2020, Hewlett Packard Enterprise installed a new HPC-System at HLRS which will be 3.5 times faster than Hazel Hen when fully deployed. This supercomputer, called Hawk, will be the fastest in the world for industrial production, powering computational engineering and research across science and industrial fields to advance applications in energy, climate, mobility, and health. The 5,632-node Hawk system will have a theoretical peak performance of about 26 Petaflops.

The deployment of the workflow as shown in this document was done on Hazel Hen. However, it can be transferred to Hawk.

As a mid-term permanent storage for the user data, the Quobyte<sup>5</sup> system is used. It is a parallel file system storage for HPC. The Quobyte software uses standard servers, disks and SSDs to create a single-point-of-failure free distributed network-based file system. During the Open Forecast project, the raw capacity is 4.8PB. The Quobyte system can either be accessed by a Quobyte native FUSE (Filesystem in User SpacE) client, by NFS, CIFS or by S3. The Quobyte system can be accessed from the HLRS HPC systems, typically by using the native client. Each Quobyte server has been connected to the Hazel Hen core network using two 25Gb Ethernet links.

## 3.3 Difference to classical computing

The work with a supercomputer is significantly different from the work with a personal computer. Instead of a graphical user interface (GUI), a Unix/Linux console is used. In addition, Hazel Hen has not been an ordinary cluster. It belonged to the Cray XC40 product family and has had an original architecture. Hazel Hen has had special nodes of the several types:

- external login nodes;
- service nodes;
- compute nodes.

External login nodes intended for getting access to the Cray XC40 system (c.f. Figure 1). On the login nodes, the user compiles code and submits jobs to the batch scheduler. To run a job on Hazel Hen, the user submits a script (or interactive job) to the batch scheduler by using the *qsub* command. It is the PBS torque submission command for batch job scripts. However, this script (or shell in an interactive job) does not run directly on the compute nodes – it first runs on a service node. "Service nodes" is a general term for non-compute nodes. The service nodes which launch jobs are more specifically called "non-Cray machine-oriented miniserver (MOM)" nodes. Like the login nodes, service nodes are not the compute nodes that make up the main computational resources of the machine but are an intermediate node where the submission script launches executables to the compute nodes with the *aprun* command. *aprun* is the ALPS (Application Level Placement Scheduler) application launcher.

<sup>&</sup>lt;sup>5</sup> https://www.quobyte.com/high-performance-computing-storage



Cray XC40



Figure 1: Architecture of the Cray XC40 system.

The only way to start a parallel job on the compute nodes is to use the portable batch system (PBS). The installed batch system on Hazel Hen was based on:

- the resource management system torque and
- the scheduler *moab*.

Production jobs are typically run in batch mode. Batch scripts are shell scripts containing flags and commands to be interpreted by a shell and are used to run a set of commands in sequence. The number of required nodes, cores, wall time and more can be determined by the parameters in the job script header with "#PBS" before any executable commands in the script.

After launch, the batch script is not necessarily granted resources immediately. It may sit in the queue of pending jobs for some time before its required resources become available. The Hazel Hen system has had different queues, e.g.:

- debugging queues (test, single);
- production queues (multi).

For example, the user can use the *test* queue for debugging of applications on up to 384 nodes, if his application finishes in less than 25 minutes. Starting the execution of an application on a large number of compute nodes can take a lot of time, depending on the available resources, and can last from one to several days. Therefore, computations on demand are not possible. Applications are run according to the schedule of the PBS.

Name	Max. #CPUs/ nodes	Max. #nodes	Max. elapsed time	Mode
test	24	384	25 minutes (CPU)	shared
single	24	1	24 hours	shared
multi	24	4096	24 hours	dedicated



An important feature of computing on supercomputers is the storage of user data. Usually users store their data on a special disk storage. The Hazel Hen has had:

- HOME directories on a shared RAID system with a small quota;
- SCRATCH directories for large files and fast I/O on a Lustre file system.

Scratch directories are available on all compute and login nodes via the workspace mechanism. This mechanism allows you to keep data outside your home not only during a simulation, but also afterwards. The idea is to allocate disk space for a number of days, and giving it a name, which allows you to identify a workspace, and to distinguish several workspaces. It is also possible, on special request only, to allocate workspaces on different file systems, which are prepared for workspaces on the local host. But workspaces have some restrictions: First, there is a maximum time limit for each workspace (60 days) after which they will be deleted automatically. Second, they have a group quota limit for space and number of files, e.g. one million files by default. Based on these restrictions, the Quobyte system as a mid-term storage was selected.

Accordingly, all data for the PMFS use-case is now stored on the Quobyte system. Security policy of HLRS does not allow any web-services to directly access this data. In addition, copying large amounts of data (terabytes) directly to a web server is also not possible: in this case, the web server needs to have its own storage and will spend a lot of time on the copying process. Therefore, it was decided to store all the data on the Quobyte system and refuse to use web services for the time being.

## 4 Use Case I: PMFS – Particulate Matter Forecast Service

The aim of the project is to develop and set up a prototype for particulate matter forecasts over the Stuttgart Metropolitan area. The targeted resolution is the turbulence permitting scale, i.e. a horizontal resolution of 50-100 m. The size of the applied domains is determined by the idea to potentially apply the PMFS as an operational forecasting system. As the computational cost of running a Large-Eddy-Simulation (LES) domain including atmospheric chemistry is very high, the domain sizes are kept as small as possible but as large as necessary. The model domains for the PMFS have been set up and the WRF-Chem model (*Grell et al., 2005; Skamarock et al., 2019*) was successfully compiled at HLRS.

Due to the target resolution of 50 m over the Stuttgart metropolitan area, we apply two outer coarser domains with 250 m and 1250 m resolution. The grid dimensions are 800x800, 601x601 and 601x601 cells. Figure 2 illustrates our domain setup.



Figure 2: Domain setup of the PMFS over the Stuttgart metropolitan area. Upper row: 1250 m resolution (left) and 250 m resolution (right). Lower panel: 50 m resolution.

The communication between the nests is one-way, that is information from domain 1 is passed as lateral boundary condition to domain 2 and information from domain 2 is passed as lateral boundary condition to domain 3.. There is no feedback from the inner domains to the corresponding parent domain.

The vertical resolution is 100 levels up to 50 hPa (22 km above sea level) with a high density in the lowest 2000 m above the surface. Meteorological input data are obtained from the operational ECMWF<sup>6</sup> analysis available on 137 vertical levels. Using this data for research applications requires an agreement with ECMWF and DWD<sup>7</sup> to get access to the ECMWF Meteorological Archival and Retrieval System (MARS) system. In case this data will be applied for real time applications, a fee-based special agreement between the institution, DWD, and ECMWF is required. To further speed up the simulation, adaptive time stepping is applied. The maximum allowed time step for the outermost model domain is 15 s while it is limited to 3 s for the innermost model domain.

<sup>&</sup>lt;sup>6</sup> https://www.ecmwf.int

<sup>7</sup> https://www.dwd.de



As the targeted horizontal resolution is very high, the so far available data set for terrain and land cover do not have a sufficient horizontal resolution. Therefore, land cover data for all domains is obtained from the Copernicus CLC 2012<sup>8</sup> data set which is aggregated to a grid with 100 m resolution. Depending on the applied resolution, terrain information is obtained from the Global Multi-resolution Terrain Elevation Data (GMTED2010<sup>9</sup>) data set as well as from the Shuttle Radar Topography Mission (SRTM<sup>10</sup>). Additionally, we incorporated a new land-cover data set from Landesamt für Umwelt Baden-Württemberg (LUBW<sup>11</sup>), which is derived from LANDSAT<sup>12</sup> in 2010 and is available at 30 m resolution suitable for our target domain. The original LUBW categories were mapped to the IGBP-MODIS<sup>13</sup> categories required by the WRF model using the QGIS software package. Figure 3 shows the land cover of the innermost model domain focusing on the Stuttgart Metropolitan area.

<sup>&</sup>lt;sup>8</sup> https://land.copernicus.eu/pan-european/corine-land-cover/clc-2012

<sup>&</sup>lt;sup>9</sup> https://www.usgs.gov/land-resources/eros/coastal-changes-and-impacts/gmted2010

<sup>10</sup> https://www2.jpl.nasa.gov/srtm/

<sup>&</sup>lt;sup>11</sup> https://www.lubw.baden-wuerttemberg.de/

<sup>12</sup> https://landsat.gsfc.nasa.gov

<sup>13</sup> https://icdc.cen.uni-hamburg.de/daten/land/modis-landsurfacetype.html



Figure 3: LUBW land cover data set prepared for the WRF simulations.

Based on the long-term experience of the Institute of Physics and Meteorology (IPM) of the University of Hohenheim with WRF in different regions of Europe and the world, we decided to apply the following physics packages.

Shortwave and longwave radiation is parameterized by the RRTMG scheme (*lacono*, 2008), which interacts with the cloud microphysics and aerosols. The Noah-MP land-surface model (LSM) (*Niu* et al., 2011) calculates soil and surface fluxes as well as soil temperatures and soil moisture coupled to the surface layer parametrization. The outermost model domain applies the YSU planetary boundary layer (PBL) parametrization (*Hong* et al., 2006), while the other two domains resolve turbulence directly without any PBL scheme. Cloud microphysics is simulated by the Thompson 2-moment cloud microphysics (*Thompson* et al., 2008) propagating hydrometeors and its number concentrations. No cumulus parametrization is applied, as the model is able to explicitly resolve convection.

The Regional Acid Deposition Model, 2nd generation (RADM2), parameterizes the atmospheric chemistry in all domains. RADM2 features 21 inorganic and 42 chemical species including more than 100 chemical reactions. Aerosols are represented by the Modal Aerosol Dynamics Model for Europe (MADE) and Secondary Organic Aerosol Model (SORGAM) scheme (chem\_opt = 106) considering size distributions, nucleation, coagulation, and condensational growth. The scheme is called every 2 min for the outer domain and every 45 s in the inner domains. Due to the applied



solver, this option allows for a larger time step for the chemistry module compared to other even more simple schemes.

Recent results indicate that the application of an urban canopy model (*Tewari* et al., 2004) improves the forecast quality especially in terms of surface fluxes. This requires special data categories for low and high residential areas as well as for industrial areas that are available in the CLC 2012 and LUBW data sets. As the finest resolution applied for the PMFS is 50 m, the Building Effect Parameterization (BEP) is not applied. The resolution is still too coarse in order to be able to resolve individual buildings.

In order to investigate the performance of the model system, scaling tests have been performed on the XC 40 system at HLRS. The left panel of Figure 4displays an example scaling of the WRF-CHEM model with respect to 20 nodes for the smallest, innermost domain. Less than 20 nodes are not possible, as the memory requirements per core would exceed the available memory of 128 GB /node.



Figure 4: Speedup of the WRF-CHEM simulation w.r.t. to 480 cores (left) and memory consumption/core (right).

It is visible that the scaling is not linear and it appears that MPI communication becomes the driving factor when applying more than 80 nodes. Compared to a WRF simulation, the required wall time per time step in our configuration is around 4-5 times higher. Due to the high resolution and thus small model integration time steps, a 24 h forecast currently requires 40 hours wall clock time. The total data amount for a 24 h forecast using 5 min output intervals including 40 chemical constituents is about five TB for the innermost model domain. Currently, the hybrid mechanism of MPI and OPENMP cannot be used for WRF-Chem but it is worthwhile to tackle this issue together with the HLRS in a future collaboration. This could lead to a reduction of the necessary wall clock time by a factor of up to four.

To initialize the model and provide chemical boundary conditions, data from the Whole Atmosphere Community Climate Model (WACCM<sup>14</sup>) are applied using the MOZBC conversion tool. Recently, high-resolution emission data for Europe from the Copernicus Atmosphere Monitoring Service (CAMS<sup>15</sup>) became available. The resolution is approx. 7x7 km and is based

<sup>14</sup> https://www2.acom.ucar.edu/gcm/waccm, Marsh et al. 2013

<sup>&</sup>lt;sup>15</sup> https://atmosphere.copernicus.eu/, https://eccad3.sedoo.fr/, Granier et al., 2019



on data from 2016. This product provides emissions of PM10, PM2\_5, SO<sub>2</sub>, CO, NO, NO<sub>2</sub>, and CH<sub>4</sub>. As the data cannot be ingested directly with the available procedures, the data conversion to the WRF-Chem data format is to be performed separately via the Earth System Model Framework (ESMF) interpolation utilities.

#### 4.1 PMFS Workflow

#### 4.1.1 Requirement analysis

The Particulate Matter Forecast Service (PMFS) is based on the Weather Research and Forecast (WRF) model version 4.0.3 together with its chemistry enhancement and, like many scientific applications for simulation, consists of the following components:

- WRF Pre-processing;
- WRF Solver<sup>16</sup>;
- Post processing<sup>17</sup>.

The main goal is to perform the cooperation of these components as a workflow. There are several requirements to the HPC workflow for the PMFS.

- 1. Data size: Keeping the large amount of input/output data. One PMFS case can consume up to 10 TB of disk space. But if one focuses only on the innermost model domain, this can be considerably reduced to about 2TB per 24h simulation.
- 2. Compute time: The large compute time, that is about 1.7 days on 2400 CPU cores.
- 3. License restrictions: We need to consider that running PMFS with ECMWF operational analysis forcing in general requires a special agreement between the institution conducting the simulation for research applications. In case of an operational forecasting system, additional cost will apply.

Data size and compute time exclude the use of "traditional" clusters, but require larger HPC resources like Hazel Hen at HLRS. Using a HPC system however imposes additional constraints, that were described already in section 3.2. Especially, the Hazel Hen system at HLRS also imposes security restrictions.

- 1. All the data will be stored on the Quobyte system rather than on workspaces. In addition, input data must be collected on the Quobyte before the PMFS is launched.
- The using web-services for the PMFS use case is technically impossible. Security policy
  of HLRS does not allow any web-services to directly access this data. All the data will be
  placed on the Quobyte system and each component of the PMFS workflow will be
  launched on HLRS resources.

<sup>&</sup>lt;sup>16</sup> Note that WRF Solver an abstract naming like in the UML diagram and refers to an executable simulation code that runs on Hazel Hen. The so-called solver actually consists of two parts: First a physical model, which, in our case, is the WRF forecast model for the climate part. Second the actual numerical solver, that solves the resulting differential equations.

<sup>&</sup>lt;sup>17</sup> In this case WRF Post Processing consists of only visualization of results.



3. The waiting time of a job in a queue also cannot be predicted. There are no computations "on demand".

From this background the pipeline for doing computations must be adapted to perform the required needs and to respect the peculiarities of Hazel Hen. It will be discussed in the following implementation section.

4.1.2 Implementation

The simplified workflow is given in

Figure 5. The first step is a pre-processing consisting of horizontal and vertical interpolations of meteorological and chemical data for initial and boundary conditions. Then the WRF simulation is performed followed by the visualization of the simulation data according to user requirements-



for the PMFS use case

The Unified Modelling Language (UML) deployment diagram for the PMFS use case is given in Figure 6. A more complex view on the necessary steps is shown in

Figure 7. In the first step all data should be copied from external data services to the Quobyte system. The second step is to run WRF Pre-processing Service on the Hazel Hen. This service includes calculations of a horizontal interpolation and a vertical interpolation. Note, that the geographical step needs to be performed only once to define the domains. Next steps are to launch WRF Solver and Post Processing Service. Each component is launched using a script written in the bash language. Data transfer between the WRF program components are handled in the already existing file format. Additional parameters and configuration files, if needed, are planned to be written by using JSON.





Figure 6: UML Deployment Diagram for the PMFS use case.

Individual workflow components can be implemented by different applications. In this case, the following applications were used.

WRF Pre-processing:

- UNGRIB reads data from different meteorological analyses and writes the data into an intermediate format that is processed from the next pre-processing step. Serial computation on one core, short time.
- METGRID interpolates the meteorological data extracted from *UNGRID* horizontally to the model domain created with *GEOGRID* and interpolates the wind field to the numerical grid applied by the WRF model.
- REAL.exe: Parallel computation on 200 cores, around 15min to prepare initial and boundary conditions.

WRF Solver:

• WRF-CHEM contains the WRF model with in-line chemistry. Parallel computation on 2400 cores, takes around 1.7 days for a 24h forecast due to limitations of the WRF I/O implementation HLRS is currently investigating the reasons for this issue.

WRF Visualization<sup>18</sup>:

<sup>&</sup>lt;sup>18</sup> A more detailed description of the results visualization can be found in the deliverable Open\_Forecast\_Visualisation\_M08.pdf on the Action's website https://open-forecast.eu/wp-content/uploads/2020/03/Open\_Forecast\_Visualisation\_M08.pdf



- COVISE
- VISTLE
- NCL

#### The data used in the PMFS use case has the following size:

- Geographical input data ~ few hundred MB;
- Meteorological input data ~ 1.5 GB per simulation day;
- Input data for the WRF solver ~ 35 GB;
- Reduced output data for the WRF solver (one day prediction) ~ 2 TB in NetCDF format.



Figure 7: Workflow to set up and operate the PMFS within Open Forecast.

#### 4.1.3 PMFS Case Study

The recently selected period serving as a case study is January 21, 2019. This day was dominated by cold air masses over Germany leading to very low temperatures during the night associated with light easterly winds and with high particulate matter and  $NO_x$  concentrations. Appendix A shows first results of selected variables of this case study.

Further enhancements may include higher resolution aerosol and chemical initial data from the CAMS regional air quality project<sup>19</sup>. It is also desirable to have near-real time emission data not only from road traffic, as this would help to considerably improve the air quality forecasts.

<sup>19</sup> https://www.regional.atmosphere.copernicus.eu/



## 5 Use Case II: AgriCOpen – Satellite Data Service for Agriculture

The use-case AgriCOpen develops seamless access to analysis-ready products derived from remote-sensing techniques to support further application in precision farming and farm management information systems (FMIS). These products are a) vegetation indices for each day of observation (CURVEG), b) a mosaic for entire Germany based on one flyover period (CURSCENE), c) a 14-day composite (CURBEST) and d) an average of the annual vegetation index development for measurements of long-term vegetation development (LONGVEG).

#### 5.1 Requirement analysis

The estimation for the required disk space resulted in an amount of 4.0 TB disk space per year. This includes all sentinel 2 scenes of A&B for entire Germany calculated with an average scene size of 850MB x 68 (tiles covering Germany) x 73 scenes per year in compressed form. After extraction the full disk space would amount to 6.1 TB. After pre-processing, it can be assumed that the disk space is reduced as one product of one single tile amount to approximately 180 MB. In sum per tile we estimate 80x8 = 45 GB per year; covering Germany by the use of all 68 scenes max. 3TB.

The permanence and processing time of atmospheric correction and cloud masking was tested on Dell Latitude 559 using Intel® Core<sup>™</sup> i5-8250U CPU @ 1.60GHz × 8 and 16GB Ram on 256 MB SSD. The entire process for atmospheric correction needed up to 40 minutes per scene. Index calculation is usually much faster and performed between minutes for an entire scene. The performance of cloud masking is tested along the implementation of GRASS and GDAL.

#### 5.2 Implementation

Based on the requirement analysis the project decided to use Sen2-Agri<sup>20</sup> system that is provided by the ESA. The system allows to download and process in parallel Sentinal-2 and Landsat-8 time series with low effort. The Sen2-Agri system can therefore be used to compute parts of the AgriCOpen use case. The system requirements for Sen2-Agri when using it for a national case, as intended, are:

- 15TB of disk space
- 64GB of RAM
- 2 quad-core CPUs (8 cores in total)

The system uses a PostgreSQL data base and installs its own SLURM<sup>21</sup> environment. The software is designed for CentOS<sup>22</sup> 7 only.

Sen2-Agri is currently installed and runs on a powerful virtual machine in the cloud environment. However, it is intended to transfer some tasks to the HPC-Systems at GWDG. Therefore, an

<sup>20</sup> http://www.esa-sen2agri.org

<sup>&</sup>lt;sup>21</sup> SLURM is an open source, fault-tolerant, and highly scalable cluster management and job scheduling system for large and small Linux clusters (https://Slurm.schedmd.com/)

<sup>&</sup>lt;sup>22</sup> https://www.centos.org



interface to connect the HPC-SLURM environment needs to be developed and is currently work in progress.

The remaining steps of the AgriCOpen use case will run on the HPC environment directly. Single steps of the workflow will be described in Singularity<sup>23</sup> containers. These containers can then be called and executed via the HPC queuing system.

## 6 Acknowledgements

We acknowledge use of the WRF-Chem preprocessor tool MOZBC provided by the Atmospheric Chemistry Observations and Modeling Lab (ACOM) of NCAR. We are also grateful to ECMWF for providing the operational analysis data on model levels.

The work carried out by the Open Forecast project is co-financed by the Connecting Europe Facility of the European Union under action number 2017-DE-IA-0170.

## 7 References

- Granier, C., S. Darras, H. Denier van der Gon, J. Doubalova, N. Elguindi, B. Galle, M. Gauss, M. Guevara, J.-P. Jalkanen, J. Kuenen, C. Liousse, B. Quack, D. Simpson, K. Sindelarova <u>The Copernicus Atmosphere Monitoring Service global and regional emissions (April 2019 version)</u> Report April 2019 version, null, *doi:10.24380/d0bn-kx16*, 2019
- Georg A. Grell, Steven E. Peckham, Rainer Schmitz, Stuart A. McKeen, Gregory Frost, William C. Skamarock, Brian Eder, 2005: Fully coupled "online" chemistry within the WRF model, Atmospheric Environment **39**, 6957-6975.
- Hong, S., Y. Noh, and J. Dudhia, 2006: <u>A New Vertical Diffusion Package with an Explicit</u> <u>Treatment of Entrainment Processes.</u> *Mon. Wea. Rev.*, **134**, 2318–2341, <u>https://doi.org/10.1175/MWR3199.1</u>
- Iacono, M. J., Delamere, J. S., Mlawer, E. J., Shephard, M. W., Clough, S. A., and Collins, W. D. 2008, Radiative forcing by long-lived greenhouse gases: Calculations with the AER radiative transfer models, *J. Geophys. Res.*, 113, D13103, doi:10.1029/2008JD009944.
- Marsh, D. R., M. Mills, D. Kinnison, J.-F. Lamarque, N. Calvo, and L. Polvani, 2013: Climate change from 1850 to 2005 simulated in CESM1(WACCM), J. Clim., 26(19), 7372– 7391, doi:10.1175/JCLI-D-12-00558.
- Niu, G-Y., Yang, Z. L., Mitchell, K. E., Chen, F., Ek, M. B., Barlage, M., Xia, Y. (2011). The community Noah land surface model with multiparameterization options (Noah-MP): 1. Model description and evaluation with local-scale measurements. *Journal of Geophysical Research: Space Physics*, *116*(12), [D12109]. https://doi.org/10.1029/2010JD015139
- Thompson, G., P.R. Field, R.M. Rasmussen, and W.D. Hall, 2008: <u>Explicit Forecasts of</u> <u>Winter Precipitation Using an Improved Bulk Microphysics Scheme. Part II:</u>

<sup>23</sup> https://singularity.lbl.gov



Implementation of a New Snow Parameterization. Mon. Wea. Rev., **136**, 5095–5115, https://doi.org/10.1175/2008MWR2387.1

- Tewari M, Chen F, Wang W, Dudhia J, LeMone MA, Mitchell K, Ek M, Gayno G, Wegiel J, Cuenca RH ,2004: Implementation and verification of the unified NOAH land surface model in the WRF model. In: 20th conference on weather analysis and forecasting/16th conference on numerical weather prediction, pp 11–15.
- Skamarock, W. C., J. B. Klemp, J. Dudhia, D. O. Gill, Z. Liu, J. Berner, W. Wang, J. G. Powers, M. G. Duda, D. M. Barker, and X.-Y. Huang, 2019: A Description of the Advanced Research WRF Version 4. NCAR Tech. Note NCAR/TN-556+STR, 145 pp. doi:10.5065/1dfh-6p97





## Appendix A: Results of the first case study

Figure 8: PM10 concentrations (ug/m<sup>3</sup>) and 10-m wind velocities at 06 UTC on January 21, 2019.



NO+NO2 concentration 2019-01-21\_06:00:00

Figure 9: Same as Fig. 8 but here for NOx. The NO<sub>x</sub> concentrations show high values of more the 100  $\mu$ g/m<sup>3</sup>.



SO<sub>2</sub> concentration 2019-01-21\_06:00:00 48°54'N 48°52'N 48°50'N 48°48'N 48°46'N 48°44'N AK Stgt. 48°42'N 48°40'N 9°E 9°5'E 9°10'E 9°15'E 9°20'E ug m<sup>-3</sup> 1 1.5 2 2.5 3 3.5 4 4.5 5 5.5 6 6.5 7 7.5 8 8.5 9 5 © T. Schwitalla Reference Vector

Figure 10: Same as Fig. 8 but here for SO<sub>2</sub>. The SO<sub>2</sub> concentrations appear to be accumulated in the Neckar valley.





Figure 11: Same as Fig. 8 but here for  $NH_3$ . The values are relatively low as they are mainly related to agriculture which is on winter break during January.





Figure 12: Zoom into Stuttgart downtown showing 2-m temperature and 10-m wind field. Due to the high resolution of 50 m, the lower temperature over the elevations can be clearly identified (blueish colors).



Ground heat flux @ 2019-01-21\_06:00:00

#### © T. Schwitalla

Figure 13: Ground heat flux. Positive values indicate heat transfer from the subsurface to the surface. Due to the applied high resolution of 50 m, the different behavior of vegetated, urban and industrial areas in terms of heat transfer is clearly visible.



PM10 concentration 2019-01-21\_07:30:00



2 4 6 8 10 12 14 16 18 20 22 24 26 28 30 32 34 Figure 14: Vertical cross section of PM10 at 07:30 UTC (8:30 am local time). The center latitude is Stuttgart main station. The Black regions denote the model terrain. It is clearly seen, that the PM10 concentration accumulates up to 200 m above ground while it reaches almost zero above 1000 m above sea level.