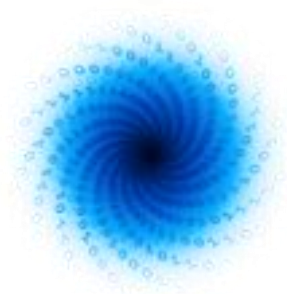




EuroHPC  
Joint Undertaking



# MAchinE Learning for Scalable meTeoROlogy and climate



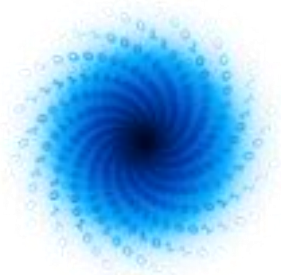
MAELSTROM

## Data Management Plan

---

Daniel Thiemert

[www.maelstrom-eurohpc.eu](http://www.maelstrom-eurohpc.eu)



MAELSTROM

## D4.2 Data Management Plan

<b>Author(s):</b>	Daniel Thiemert (ECMWF)
<b>Dissemination Level:</b>	Public
<b>Date:</b>	21/09/2021
<b>Version:</b>	1.0
<b>Contractual Delivery Date:</b>	30/09/2021
<b>Work Package/ Task:</b>	WP4/ T4.1
<b>Document Owner:</b>	ECMWF
<b>Contributors:</b>	All Partners
<b>Status:</b>	Final



# MAELSTROM

## Machine Learning for Scalable Meteorology and Climate

**Research and Innovation Action (RIA)**

**H2020-JTI-EuroHPC-2019-1: Towards Extreme Scale Technologies and Applications**

**Project Coordinator:** Dr Peter Dueben (ECMWF)

**Project Start Date:** 01/04/2021

**Project Duration:** 36 months

**Published by the MAELSTROM Consortium**

**Contact:**

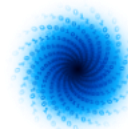
ECMWF, Shinfield Park, Reading, RG2 9AX, United Kingdom

[Peter.Dueben@ecmwf.int](mailto:Peter.Dueben@ecmwf.int)

The MAELSTROM project has received funding from the European High-Performance Computing Joint Undertaking (JU) under grant agreement No 955513. The JU receives support from the European Union's Horizon 2020 research and innovation programme and United Kingdom, Germany, Italy, Luxembourg, Switzerland, Norway



**EuroHPC**  
Joint Undertaking

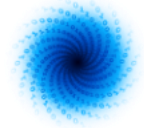


## Contents

<b>1</b>	<b>EXECUTIVE SUMMARY</b> .....	<b>5</b>
<b>2</b>	<b>INTRODUCTION</b> .....	<b>6</b>
2.1	ABOUT MAELSTROM.....	6
2.2	SCOPE OF THIS DELIVERABLE .....	6
2.2.1	<i>OBJECTIVES OF THIS DELIVERABLE</i> .....	6
2.2.2	<i>WORK PERFORMED IN THIS DELIVERABLE</i> .....	6
2.2.3	<i>DEVIATIONS AND COUNTER MEASURES</i> .....	6
<b>3</b>	<b>OPEN RESEARCH DATA OBJECTIVES</b> .....	<b>7</b>
3.1	OPEN RESEARCH DATA PILOT .....	7
3.2	MAELSTROM RESEARCH DATA.....	7
3.3	DATA MANAGEMENT PLAN QUESTIONNAIRE .....	8
<b>4</b>	<b>MAELSTROM DATA SETS</b> .....	<b>10</b>
<b>5</b>	<b>CONCLUSION</b> .....	<b>19</b>

## Tables

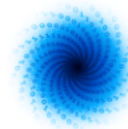
Table 1:	MAELSTROM Data Management Plan Questionnaire .....	8
----------	--	---



## 1 Executive Summary

The MAELSTROM Data Management Plan responds to the requirements of the H2020 Open Research Data Pilot to document which research data is being produced by the MAELSTROM project, in which format, and how it will be made available.

It has already identified data sets envisaged by the project but is only to be seen as an initial version which requires periodic updates to provide the necessary detail as it emerges.



## 2 Introduction

### 2.1 About MAELSTROM

To develop Europe's computer architecture of the future, MAELSTROM will co-design bespoke compute system designs for optimal application performance and energy efficiency, a software framework to optimise usability and training efficiency for machine learning at scale, and large-scale machine learning applications for the domain of weather and climate science.

The MAELSTROM compute system designs will benchmark the applications across a range of computing systems regarding energy consumption, time-to-solution, numerical precision and solution accuracy. Customised compute systems will be designed that are optimised for application needs to strengthen Europe's high-performance computing portfolio and to pull recent hardware developments, driven by general machine learning applications, toward needs of weather and climate applications.

The MAELSTROM software framework will enable scientists to apply and compare machine learning tools and libraries efficiently across a wide range of computer systems. A user interface will link application developers with compute system designers, and automated benchmarking and error detection of machine learning solutions will be performed during the development phase. Tools will be published as open source.

The MAELSTROM machine learning applications will cover all important components of the workflow of weather and climate predictions including the processing of observations, the assimilation of observations to generate initial and reference conditions, model simulations, as well as post-processing of model data and the development of forecast products. For each application, benchmark datasets with up to 10 terabytes of data will be published online for training and machine learning tool-developments at the scale of the fastest supercomputers in the world. MAELSTROM machine learning solutions will serve as blueprint for a wide range of machine learning applications on supercomputers in the future.

### 2.2 Scope of this deliverable

#### 2.2.1 Objectives of this deliverable

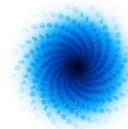
D4.2 Data Management Plan provides the initial outline of the data management plan including information on which data sets will be created in the project and how they will be made available. This document represents only the initial version where details may not be available yet, and it will be further developed over the course of the project.

#### 2.2.2 Work performed in this deliverable

The work performed included, as per the DoA, the collection of the available descriptions of data sets to be produced by the project, through a questionnaire.

#### 2.2.3 Deviations and counter measures

No deviations have been encountered.



## 3 Open Research Data Objectives

### 3.1 Open Research Data Pilot

As per the Guidelines to the Rules on Open Access to Scientific Publications and Open Access to Research Data in Horizon 2020<sup>1</sup>, Research Data

“Refers to information, in particular facts or numbers, collected to be examined and considered as a basis for reasoning, discussion, or calculation. In a research context, examples of data include statistics, results of experiments, measurements, observations resulting from fieldwork, survey results, interview recordings and images. The focus is on research data that is available in digital form.”

The Open Research Data Pilot

“aims to improve and maximise access to and re-use of research data generated by Horizon 2020 projects<sup>2</sup>”

and applies to data sets that are

“needed to validate the results presented in scientific publications<sup>2</sup>”.

The Data Management Plan is expected to

“specify what data will be open: detailing what data the project will generate, whether and how it will be exploited or made accessible for verification and re-use, and how it will be curated and preserved<sup>2</sup>”.

### 3.2 MAELSTROM Research Data

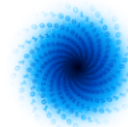
As per the Description of Actions, MAELSTROM will generate and use model output from Weather & Climate (W&C) models, weather observations (including standard products but also crowd-sourced weather station data), data from energy production (e.g., wind farms or solar panels), as well as data from citizen sensing, crowd-sensing and tweeting sensors. Only de-personalised/anonymised data will be handled within MAELSTROM. The data will be unified via a data API that is organised by WP2 and mainly based on the CliMetLab tool<sup>3</sup>. The data will be used to generate benchmark datasets of MAELSTROM that will be published as open access data via S3 buckets. The datasets are made available in different tiers where possible, including a tier-1 version that can be used for training purposes that will not exceed a size that is not manageable by a standard laptop and larger tiers that can be used for scaling experiments with tens of terabytes for some of the applications. Since the data will be available as open source, verification and re-use will be guaranteed. The datasets will be curated by ECMWF as they are meant to become a valuable asset of the W&C community to allow

---

<sup>1</sup> [http://ec.europa.eu/research/participants/data/ref/h2020/grants\\_manual/hi/oa\\_pilot/h2020-hi-oa-pilot-guide\\_en.pdf](http://ec.europa.eu/research/participants/data/ref/h2020/grants_manual/hi/oa_pilot/h2020-hi-oa-pilot-guide_en.pdf)

<sup>2</sup> [https://ec.europa.eu/research/participants/docs/h2020-funding-guide/cross-cutting-issues/open-access-dissemination\\_en.htm](https://ec.europa.eu/research/participants/docs/h2020-funding-guide/cross-cutting-issues/open-access-dissemination_en.htm)

<sup>3</sup> <https://climetlab.readthedocs.io/en/latest/>



for quantitative evaluations of Machine Learning (ML) solutions. ECMWF will cover the cost for data curation.

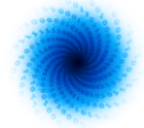
### 3.3 Data Management Plan Questionnaire

The following questionnaire has been provided to MAELSTROM technical work packages to gather the information for this first version of the Data Management Plan.

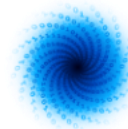
*Table 1: MAELSTROM Data Management Plan Questionnaire*

<b>&lt;Data set reference and name&gt;</b>	
<b>Data set description</b>	<p><i>Description of the data that will be generated or collected (or is already available to the project), its origin (in case it is collected), nature and scale and to whom it could be useful, and whether it underpins a scientific publication. Information on the existence (or not) of similar data and the possibilities for integration and reuse.</i></p> <p><i>Limitations?</i></p> <p><i>Constraints?</i></p>
<b>Standards and metadata</b>	<p><i>Reference to existing suitable standards of the discipline. If these do not exist, an outline on how and what metadata will be created.</i></p> <p><i>Will you generate proper metadata for you data?</i>  <i>If yes: how do they look like?</i>  <i>If no: why?</i></p> <p><i>Data format?</i></p> <p><i>Will there be a review process to quality- check the data?</i></p>
<b>Data Sharing</b>	<p><i>Description of how data will be shared, including access procedures, embargo periods (if any), outlines of technical mechanisms for dissemination and necessary software and other tools for enabling re-use, and definition of whether access will be widely open or restricted to specific groups. Identification of the repository where data will be stored, if already existing and identified, indicating in particular the type of repository (institutional, standard repository for the discipline, etc.).</i></p> <p><i>In case the dataset cannot be shared, the reasons for this should be mentioned (e.g. ethical, rules of personal data, intellectual property, commercial, privacy-related, security-related).</i></p> <p><i>License?</i></p> <p><i>Access URL?</i></p>
<b>Archiving and preservation (including storage and backup)</b>	<p><i>Description of the procedures that will be put in place for long-term preservation of the data. Indication of how long the data should be preserved, what is its approximated end volume, what the associated costs are and how these are planned to be covered.</i></p>





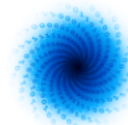
	<p><i>At which Data Center do you want to store your data?</i></p> <p><i>Is there an established workflow for your requested DOI process in place?</i></p> <p><i>According to which standards</i></p>
--	---



## 4 MAELSTROM Data Sets

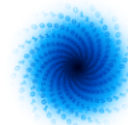
The following sections provide the responses by work packages.

<b>maelstrom-yr</b>	
<b>Data set description</b>	<p>Contains gridded weather data for the Nordics. It contains both predictors (gridded weather forecasts) and predictand (gridded analysis fields). The forecasts are used operationally for the Nordic region on <a href="https://www.yr.no">https://www.yr.no</a> .</p> <p>The use of a subset of this dataset has been described in “Nipen, T. N., I. A. Seierstad, C. Lussana, J. Kristiansen, and Ø. Hov, 2020: Adopting Citizen Observations in Operational Weather Prediction. Bull. Amer. Meteor. Soc., 101, E43–E57.”</p> <p><i>Limitations</i></p> <p>No</p> <p><i>Constraints</i></p> <p>No</p>
<b>Standards and metadata</b>	<p>The dataformat is NetCDF and the dataset follows the NetCDF/CF standard. The dataset contains metadata that describes the coordinate system, target and predictor variables, and their units. The metadata is stored in the NetCDF files.</p> <p><i>Will you generate proper metadata for your data?</i></p> <p><i>If yes: how do they look like?</i></p> <p><i>If no: why?</i></p> <p><i>Data format</i></p> <p>netCDF</p> <p><i>Will there be a review process to quality-check the data?</i></p> <p>The review process will consist of spot checks of the data performed by an independent person not involved with the data generation.</p>
<b>Data Sharing</b>	<p>Jupyter notebooks have been created to explore this dataset. CliMetLab is used to download for data loading. The Jupyter can be accessed through github and be used for public.</p> <p><i>In case the dataset cannot be shared, the reasons for this should be mentioned (e.g. ethical, rules of personal data, intellectual property, commercial, privacy-related, security-related).</i></p> <p><i>License</i></p>



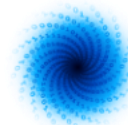
	<p>The license of the data is described here: <a href="https://www.met.no/en/free-meteorological-data/Licensing-and-crediting">https://www.met.no/en/free-meteorological-data/Licensing-and-crediting</a></p> <p><i>Access</i></p> <p><a href="https://github.com/metno/maelstrom-yr/blob/main/notebooks/demo_yr.ipynb">https://github.com/metno/maelstrom-yr/blob/main/notebooks/demo_yr.ipynb</a></p>
<b>Archiving and preservation (including storage and backup)</b>	<p><i>At which Data Center do you want to store your data?</i></p> <p>ECMWF data cloud</p> <p><i>Is there an established workflow for your requested DOI process in place?</i></p> <p><i>According to which standards</i></p> <p>Not yet.</p>

<b>maelstrom-radiation</b>	
<b>Data set description</b>	<p>This is a dataset to train machine learning emulators for parametrisation schemes of the IFS weather forecast model. The data is collected from forecast simulations with the IFS. The data has a size of several terabytes and will be useful to researchers and developers of weather and climate models. Emulators that are trained from the Deliverable can only be as good as the original scheme, unless the data is used as inputs to generate more sophisticated output fields.</p> <p>This data is for learning the emulation of the ECMWF radiation scheme, TripleClouds, found in the ecRad package (<a href="https://github.com/ecmwf/ecrad">https://github.com/ecmwf/ecrad</a>). Building an accurate emulator of radiative heating could accelerate weather and climate models partially by enabling the use of GPU hardware within our models.</p> <p><i>Limitations</i></p> <p>No</p> <p><i>Constraints</i></p> <p>No</p>
<b>Standards and metadata</b>	<p>Within the netCDF files the data follows CF conventions. For speed of learning this metadata is removed in the packing in TFRecords format, however this data is a repackaging of the netCDF which can be consulted for referencing metadata.</p> <p><i>Will you generate proper metadata for you data?</i></p> <p><i>If yes: how do they look like?</i></p> <p><i>If no: why?</i></p> <p><i>Data format?</i></p> <p>netCDF and TFRecords</p>



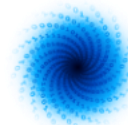
	<p><i>Will there be a review process to quality- check the data?</i></p> <p>Yes</p>
<b>Data Sharing</b>	<p>Jupyter notebooks have been created to explore this dataset. CliMetLab is used to download for data loading. The Jupyter can be accessed through github and be used for the public.</p> <p><i>License</i></p> <p><a href="https://storage.ecmwf.europeanweather.cloud/MAELSTROM_AP3/LICENCE.txt">https://storage.ecmwf.europeanweather.cloud/MAELSTROM_AP3/LICENCE.txt</a></p> <p><i>Access URL</i></p> <p><a href="https://git.ecmwf.int/projects/MLFET/repos/maelstrom-radiation/browse/notebooks/demo_radiation.ipynb">https://git.ecmwf.int/projects/MLFET/repos/maelstrom-radiation/browse/notebooks/demo_radiation.ipynb</a></p>
<b>Archiving and preservation (including storage and backup)</b>	<p>The approximated end volume will be ~ 10 TB</p> <p><i>At which Data Center do you want to store your data?</i></p> <p>ECMWF cloud storage</p> <p><i>Is there an established workflow for your requested DOI process in place?</i></p> <p><i>According to which standards</i></p> <p>Not yet, but several options to obtain a DOI are investigated.</p>

<b>maelstrom-nogwd</b>	
<b>Data set description</b>	<p>Contains the input/output dataset for learning non-orographic gravity wave drag, as described in <a href="https://agupubs.onlinelibrary.wiley.com/doi/full/10.1029/2021MS002477">https://agupubs.onlinelibrary.wiley.com/doi/full/10.1029/2021MS002477</a>. Data is grouped by forecast start date.</p> <p>Data has been pre-processed into inputs "x" and outputs "y". "x" contains vertical profiles of winds &amp; temperature plus surface values of pressure and geopotential. "y" contains the wind increments due to parametrised non-orographic gravity wave drag. The machine learning task is to predict y given x. Unlike many ML tasks within the field of weather and climate, this task can be predicted independently for each column of the atmosphere.</p> <p><i>Limitations</i></p> <p>No</p> <p><i>Constraints</i></p> <p>No</p>

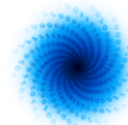


<p><b>Standards and metadata</b></p>	<p><i>Will you generate proper metadata for you data?</i></p> <p><i>If yes: how do they look like?</i></p> <p><i>If no: why?</i></p> <p>No, because the data is only available in a friendly format for machine learning, with the variables concatenated together for reading speed. This makes adding metadata challenging. A full description of the variables and their ordering is provided in the data documentation.</p> <p><i>Data format</i></p> <p>netCDF</p> <p><i>Will there be a review process to quality- check the data?</i></p> <p>Yes</p>
<p><b>Data Sharing</b></p>	<p><i>License</i></p> <p><a href="https://storage.ecmwf.europeanweather.cloud/NOGWD_L91_PUBLIC/LICENCE.txt">https://storage.ecmwf.europeanweather.cloud/NOGWD_L91_PUBLIC/LICENCE.txt</a></p> <p><i>Access URL</i></p> <p><a href="https://git.ecmwf.int/projects/MLFET/repos/maelstrom-nogwd/browse/notebooks/demo_nogwd.ipynb">https://git.ecmwf.int/projects/MLFET/repos/maelstrom-nogwd/browse/notebooks/demo_nogwd.ipynb</a></p>
<p><b>Archiving and preservation (including storage and backup)</b></p>	<p>The approximated end volume will be ~3Gb</p> <p><i>At which Data Center do you want to store your data?</i></p> <p>ECMWF cloud storage</p> <p><i>Is there an established workflow for your requested DOI process in place?</i></p> <p><i>According to which standards</i></p> <p>Not yet</p>

**maelstrom-downscaling**

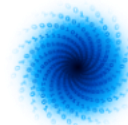


<p><b>Data set description</b></p>	<p>The two datasets will comprise various meteorological quantities which are considered to convey relevant information for downscaling the 2m temperature. In addition to the 2m temperature itself, these are temporally variable quantities such as the surface enthalpy fluxes, the cloud cover, or the temperature at different pressure levels, but also constant features such as the surface elevation.</p> <p>The data will be obtained from historical IFS HRES-forecasts (from 2016 onwards) in the short-range (max. lead time 24 hours). This provides data is used at 0.1 resolution on a spherical grid. Central Europe is chosen as the target region since the topography is highly complex in this area.</p> <p>As part of deliverable D1.1., a subset of the data (2m temperature and surface elevation at initial times of the IFS model) are used and artificially coarsened onto a 0.8° grid to imitate the input data for a supervised downscaling application.</p> <p>In the future, the application will approach a real-case application by e.g. applying ERA5 reanalysis data (on a coarser 0.3°-grid) as input data. Additionally, more quantities will reside in the dataset.</p> <p>In the future, other model data sources will be considered such as model forecast from the multiscale ICON-model operated at DWD (German Weather Service). The latter would enable supervised training from a coarse-grained global grid at 13 km resolution to a high-resolved grid at 2.2 km resolution over Germany.</p> <p><i>Limitations</i></p> <p>No</p> <p><i>Constraints</i></p> <p>No</p>
<p><b>Standards and metadata</b></p>	<p><i>Will you generate proper metadata for your data?</i></p> <p><i>If yes: how do they look like?</i></p> <p>Model data are commonly shipped along with metadata that typically follows the so-called CF-conventions. Thus, information on the data source, the underlying grid, the validity time (initial time of the model and the lead time) etc. are already part of the processed data and no further metadata creation process is required. However, pre-processing operations using CDO will also be tracked.</p> <p><i>If no: why?</i></p> <p><i>Data format</i></p> <p>netCDF</p> <p><i>Will there be a review process to quality- check the data?</i></p> <p>Yes</p>



<p><b>Data Sharing</b></p>	<p>Data sharing of the IFS HRES data is possible in scope of the project. Namely, deliverable D1.1. already realized a publication of some IFS HRES data under the Apache 2.0 License. The IFS HRES data itself was processed using the Meteorological Archival and Retrieval System (MARS).</p> <p>The ERA5 reanalysis data can be obtained from the C3S climate data store. The associated license can be viewed here:  <a href="https://apps.ecmwf.int/datasets/licences/copernicus/">https://apps.ecmwf.int/datasets/licences/copernicus/</a></p> <p>The ICON model data is distributed via DWD's Open Data Server. Terms and conditions are claimed here.</p> <p><i>License</i></p> <p>Yes</p> <p><i>Access URL</i></p> <p><a href="https://git.ecmwf.int/projects/MLFET/repos/maelstrom-downscaling-ap5/browse">https://git.ecmwf.int/projects/MLFET/repos/maelstrom-downscaling-ap5/browse</a></p>
<p><b>Archiving and preservation (including storage and backup)</b></p>	<p><i>At which Data Center do you want to store your data?</i></p> <p>The data is stored in the ECMWF cloud storage.</p> <p><i>Is there an established workflow for your requested DOI process in place?</i></p> <p><i>According to which standards</i></p> <p>Not yet.</p>

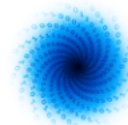
<p><b>maelstrom-power-production</b></p>	
<p><b>Data set description</b></p>	<p>A dataset plugin for CliMetLab for the dataset climetlab-plugin-a6/maelstrom-production-forecasts. It consisted of Weather forecast data and power production data and described in Deliverable D1.1.</p> <p><i>Limitations</i></p> <p>No</p> <p><i>Constraints</i></p> <p>No</p>



<p><b>Standards and metadata</b></p>	<p><i>Will you generate proper metadata for you data?</i></p> <p>Not needed, all metadata are included.</p> <p><i>Data format</i></p> <p>netCDF</p> <p><i>Will there be a review process to quality- check the data?</i></p> <p>Yes</p>
<p><b>Data Sharing</b></p>	<p>Jupyter notebooks have been created to explore this dataset. CliMetLab is used for data loading. The Jupyter can be accessed through github and be used for public.</p> <p><i>License</i></p> <p>Apache 2.0</p> <p><i>Access URL</i></p> <p><a href="https://github.com/4castRenewables/climetlab-plugin-a6/blob/main/notebooks/demo_maelstrom_production_forecasts.ipynb">https://github.com/4castRenewables/climetlab-plugin-a6/blob/main/notebooks/demo_maelstrom_production_forecasts.ipynb</a></p>
<p><b>Archiving and preservation (including storage and backup)</b></p>	<p>The approximated volume will be ~ GB - ~TB</p> <p><i>At which Data Center do you want to store your data?</i></p> <p>ECMWF cloud storage</p> <p><i>Is there an established workflow for your requested DOI process in place?</i></p> <p><i>According to which standards</i></p> <p>Not yet.</p>

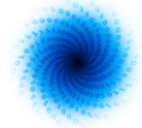
<p><b>Social media data</b></p>	
<p><b>Data set description</b></p>	<p>The data are obtained by the Twitter API. At the moment, the access is not granted, application for data access are still ongoing at Twitter.</p>
<p><b>Standards and metadata</b></p>	<p>The data are stored as obtained from Twitter through the API. We will eventually store the data in a database and dump the data as csv to the ECMWF S3 bucket.</p> <p><i>Will you generate proper metadata for you data?</i></p> <p><i>If yes:</i> we will use the incoming policy from twitter, e.g.location information, if exists. Further details are not yet clear</p>



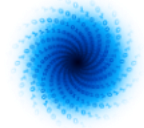


	<p><i>Data format?</i> Csv</p> <p><i>Will there be a review process to quality- check the data?</i> No</p>
<b>Data Sharing</b>	<p>We will build a database and dump data to the S3 bucket. Sharing is not allowed due to the Twitter policies. Only if Twitter has agreed, we can share.</p> <p><i>License?</i> None</p> <p><i>Access URL?</i> None</p>
<b>Archiving and preservation (including storage and backup)</b>	<p>We follow the Twitter rules, cf. <a href="https://developer.twitter.com/en/developer-terms/commercial-terms">https://developer.twitter.com/en/developer-terms/commercial-terms</a></p>

<b>maelstrom-ens10</b>	
<b>Data set description</b>	<p>ENS10 is an open dataset developed as part of MAELSTROM to enable research into machine learning of weather ensemble prediction systems via post-processing. It consists of the model output data of ECMWF "hindcast" experiments. These are global ensemble forecasts with 10 ensemble members that are spread over 20 years (1998-2017) with two forecasts per week. The dataset holds 10 ensemble members on 0.5 degree latitude/longitude grids. ENS10 weighs 3TB, thus for smaller-scale experiments we provide a smaller, tier-1 version called ENS5-mini, which weighs 9GB and consists of 5 members, 10 years, one day into the future, and is cropped over central Europe.</p> <p>One limitation of the dataset is that it does not include the full, 51-member ensemble used in production systems, nor all fields used for weather prediction. Extending ENS10 will be researched as part of this project.</p>
<b>Standards and metadata</b>	<p>The metadata of ENS10 is included within its data format – GRIB – which contains spatial information and units of each of the provided fields.</p> <p>Since the data consists of hindcasts from the ECMWF system, the information contained with it is sound. Therefore, our review process of the data consists of data-scientific inspection of the dataset and testing it on machine learning models.</p>
<b>Data Sharing</b>	<p>The dataset, in its raw GRIB format, is already uploaded onto the European Weather Cloud S3 servers and mirrored on a public ETH Zurich storage server. Both links are public and available as part of the ECMWF website, ETH Zurich, and the CliMetLab plugin we provide to accompany the dataset. The dataset is licenced under the ECMWF general data access licence (<a href="https://apps.ecmwf.int/datasets/licences/general/">https://apps.ecmwf.int/datasets/licences/general/</a>, CC BY 4.0).</p> <p>Access links:</p> <ul style="list-style-type: none"> <li>• Via CliMetLab: <a href="https://climetlab-maelstrom-ens10">climetlab-maelstrom-ens10</a></li> </ul>



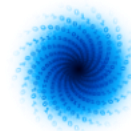
	<ul style="list-style-type: none"><li>• Via the European Weather Cloud: <a href="https://storage.ecmwf.europeanweather.cloud/MAELSTROM_AP4">https://storage.ecmwf.europeanweather.cloud/MAELSTROM_AP4</a></li><li>• Via the ETH Zurich mirror: <a href="http://spclstorage.inf.ethz.ch/projects/deep-weather/ENS10/">http://spclstorage.inf.ethz.ch/projects/deep-weather/ENS10/</a> , <a href="http://spclstorage.inf.ethz.ch/projects/deep-weather/ENS5mini/">http://spclstorage.inf.ethz.ch/projects/deep-weather/ENS5mini/</a></li></ul>
<b>Archiving and preservation (including storage and backup)</b>	<p>There is not yet a DOI for the dataset until its finalization w.r.t. fields and ensemble member counts.</p> <p>ENS10 will follow the European Weather Cloud data center policy for long-term retention. The mirror at ETH Zurich is kept as part of the Scalable Parallel Computing Laboratory servers, which contain more datasets, and will stay on for as long as the lab is operational, with no additional storage costs.</p>



## 5 Conclusion

This initial Data Management Plan has identified a number of data sets for each of the technical work packages, identifying the required details (where possible) on what data will be open, how it will be made accessible, and how it will be curated.

The Data Management Plan is to be seen as a living document and will be reviewed and revised periodically to ensure that information contained therein is up-to-date and correct.



## Document History

Version	Author(s)	Date	Changes
<b>0.1</b>	Daniel Thiemert (ECMWF)	20/09/2021	Initial version
<b>1.0</b>	Daniel Thiemert (ECMWF)	21/09/2021	Final Version

## Internal Review History

Internal Reviewers	Date	Comments
<b>Peter Dueben (ECMWF)</b>	21/09/2022	Approved with comments

## Estimated Effort Contribution per Partner

Partner	Effort
<b>ECMWF</b>	0.2
<b>Total</b>	<b>0.2</b>

This publication reflects the views only of the author, and the European High-Performance Computing Joint Undertaking or Commission cannot be held responsible for any use which may be made of the information contained therein.