



## D5.2 Data Management Plan Version 1.1

### Document Information

<b>Contract Number</b>	828947
<b>Project Website</b>	<a href="http://www.enerxico-project.eu">www.enerxico-project.eu</a>
<b>Contractual Deadline</b>	M6
<b>Dissemination Level</b>	Public
<b>Nature</b>	Report
<b>Author</b>	Josep de la Puente (BSC), Marta Rosello (BSC), Claudia Rosas (BSC)
<b>Contributor(s)</b>	All partners
<b>Reviewer</b>	
<b>Keywords</b>	data description, data collection, FAIR



**Notice:**

*The research leading to these results has received funding from the European Union's Horizon 2020 Programme under the ENERXICO Project ([www.enerxico-project.eu](http://www.enerxico-project.eu)), grant agreement no 828947 and under the Mexican CONACYT-SENER-Hidrocarburos grant agreement B-S-69926.*

## Change Log

Version	Author	Description of Change
1.0	Claudia Rosas	Initial draft
1.1	Marta Rosselló	Final draft

## Table of Contents

Executive Summary	4
<b>Introduction</b>	4
<b>Data Summary and Structure of the ENERXICO DMP Annex</b>	4
Dataset sheet	7
Software sheet	7
<b>FAIR data</b>	8
Making ENERXICO data Findable	9
Making ENERXICO data openly Accessible	10
Making ENERXICO data Interoperable	10
Increase ENERXICO data Re-use	10
<b>Allocation of Resource</b>	11
<b>Data security</b>	11
<b>Ethical aspects</b>	11
<b>Annex I</b>	11

## Executive Summary

This deliverable presents the data management plan (DMP) of the ENERXICO project, which describes the data management life-cycle for all codes and datasets to be collected, processed and/or generated during the lifetime of the project.

The document describes the contents and organization of the DMP, considering that the actual data (and metadata) description will be furnished as a “live” Annex that will be periodically updated during the development of the project. Specifically, this deliverable describes:

- How the different types of datasets that will be generated, collected, and processed during the project will be managed during and after it. This affects mainly the ENERXICO scientific codes and, to a lesser extent, other research activities.
- Which methodologies and which standards (if any) will be applied to manage each of the ENERXICO datasets.
- How the datasets will be stored and handled during the lifetime of the project and after its conclusion, as well as how the datasets will be made (openly) accessible.

## 1. Introduction

The ENERXICO project is part of the H2020 FETHPC Open Research Data Pilot (ORD pilot), aimed at improving and maximizing access and re-use of research data, while taking into account the need to balance openness and protection of scientific information, commercialization and Intellectual Property Rights (IPR), privacy concerns, security, and related data management and preservation questions.

This DMP describes how data (in its broadest sense) will be managed during the project, but it is not a final closed document. The DMP described in this report (D5.6) will be updated over the course of the project to account for the generation/acquisition of new datasets, implementation of consortium policies, and/or other external factors.

This document follows the Horizon 2020 “FAIR” DMP template and follows the FAIR data guiding principles; i.e. that data must be Findable, Accessible, Interoperable, and Re-usable.

## 2. Data Summary and Structure of the ENERXICO DMP Annex

As stated in the Grant Agreement (GA), ENERXICO will promote optimization and performance and efficiency enhancements in European and Mexican parallel codes, related to the fields of Geophysics, Combustion and Wind simulation, for the upcoming pre-Exascale and Exascale supercomputers, and will further improve the solution of scientific problems requiring of Exascale computing. Some of the final demonstrators can be considered as small-scale demonstrators for optimizing and testing codes on

Exascale hardware prototypes and for addressing the Exascale challenges.

The reference parallel applications will be the main sources of project data, which will include:

- The source code of the parallel applications themselves.
- Datasets generated to evaluate code performance.
- Experimental datasets to be used for the performance audits.
- Datasets generated from the execution of models.

For this reason, the ENERXICO DMP has been organized around a (dynamic) spreadsheet annex document, with one dedicated sheet for each of the parallel scientific codes (i.e. 8 sheets in total).

## 2.1 Dataset sheet

For each dataset in a parallel scientific code, a dataset sheet exists with information on the following data items:

Item	Comments/explanation
Name	Descriptive name to identify the dataset
Description	Short description of the contents
Data category	Data category code (see Table 2 for the corresponding codes)
License	Chosen among the most appropriated and most open ones
Repository location	Institutional or public repository name
Author	Data author(s) name(s)
Naming Conventions	File names structure and conventions
Versioning	How and where the version of the dataset can be found
Format	Standards, definitions, ontologies, etc.
Size	Total estimated size, or single file size and number of expected files
Storage	Physical support selected, dependent on availability needs
Archive path	Folders structure
Associated metadata	Selected metadata standards, and to metadata set
Provenance	Structured dataset origin information

Backup needs	Periodicity, subsets backup needs analysis, etc.
Access permissions	Lifecycle dependency: selected groups, or public
Legal/ethical restrictions	Privacy and security issues
Reproducibility	If yes: connection to code and environment
Data transfer needs	Replicas and periodic transfers to/from other repositories
Long term preservation	Needs at 3-5-7-10 years (if any)
Metadata management	Way to access metadata when data are not available
Resources need	Analysis of resources needs at each step of data lifecycle

Table 1. List of items in a dataset sheet and their definition.

Data category	Code	Name	Comments
Scientific data	1.1	Models	
	1.2	Experimental	Data coming from observation, measurements or produced by detectors/sensors or by any other experimental device and or activity.
	1.3	Synthetic	Data generated by a simulation and/or are not obtained by direct measurement
	1.4	Test	Datasets (experimental or synthetic) used to validate models
Software	2.1	Libraries	
	2.2	Applications	
	2.3	Services	
	2.4	API	
Administrative	3.1	Documents	Any documentation, either public

docs			or private, such as code documentation, technical notes, etc., not directly mentioned in the project deliverables list.
	3.2	Internal reports	Meetings minutes, internal notes to document the evolution of the project, such as calendar, resources management, mailing lists, etc.
	3.3	Deliverables	Project output documents
Other	4.1	Metadata	Any data describing data properties. If they contain scientific information, they can also be classified as scientific data.

Table 2. Summary of the different data categories

## 2.2 Software sheet

Similar to the dataset sheets, one software sheet per code (or related software component) will contain information on the following aspects:

Item	Comments/explanation
Reference name of the program or workflow	Name of the code
Description	Brief description
Author	Main developers
Programming language	Specify
Rules and best coding practices	Conventions for filenames, link to an external manual (if exists), e.g. PEP8.
Access permissions and license	Lifecycle dependency: only specific groups of collaborators, all partners, whole community, etc.
Code size if relevant (to be updated)	
Repository type	GitHub, GitLab, Bitbucket, SourceForge
Repository structure	Branches, tags, etc.
Provenance information	Containers, virtual environments

Backup and Archiving needs	If any
Legal/ethical restrictions	If any
Versioning control and rules/workflows managing	Specify the repository
Code transfer needs and security	If any
Long term preservation needs	Only if applies to a given official release version
Documentation and inline comments rules	Specify
Metadata management	(available even when the software is not)
Resources need	At each step of the life cycle

Table 3. List of items in a software sheet.

### 3. FAIR data

The FAIR Guiding Principles (Wilkinson et al.; 2016: DOI: 10.1038/sdata.2016.18) describes distinct considerations for contemporary data publishing environments with respect to supporting both manual and automated deposition, exploration, sharing, and reuse. A metric to quantify the degree of “FAIRness” of each dataset in ENERXICO has been defined. It results on a normalized value (between 0 and 1) for each of the 4 FAIR components. In turn, this (0,1) value results from assigning a flag value 0/1 to each of the FAIR subcomponents defined by Wilkinson et al. (2016) and listed in Table 4.

<b>F</b>	<b>FINDABLE</b>		
F.1	Persistent Identifiers (PDI)	(meta)data are assigned a globally unique and persistent identifier	0/1
F.2	Rich metadata	data are described with rich metadata (defined by subcomponent R.1 below)	0/1
F.3	Data registered in searchable resources	(meta)data are registered or indexed in a searchable resource	0/1
F.4	Metadata specifies the PDI	metadata clearly and explicitly include the identifier of the data it describes	0/1
<b>A</b>	<b>ACCESSIBLE</b>		

A.1	Retrievable by the PDI with a standardized protocol	(meta)data are retrievable by their identifier using a standardized communications protocol	0/1
A.1.1	Open, free protocol	The protocol is open, free, and universally implementable	0/1
A.1.2	Authentication and Authorization	The protocol allows for an authentication and authorization procedure, where necessary	0/1
A.2	Metadata availability	Metadata are accessible beyond the data availability	0/1
<b>I</b>	<b>INTEROPERABLE</b>		
I.1	Formal, accessible, shared and applicable language	(meta)data use a formal, accessible, shared, and broadly applicable language for knowledge representation	0/1
I.2	FAIR vocabulary	(meta)data use vocabulary that follow FAIR principles	0/1
I.3	Metadata references	Metadata includes qualified references to other metadata	0/1
<b>R</b>	<b>REUSABLE</b>		
R.1	Relevant metadata	(meta)data have plurality of accurate and relevant attributes	0/1
R.1.1	Usage license	(meta)data are released with clear and accessible data usage license	0/1
R.1.2	Provenance	(meta)data are associated with detailed provenance	0/1
R.1.3	Community standards	(meta)data meet domain-relevant community standards	0/1

Table 4. Definition of the different FAIR components and flag value (0/1) used to quantify the degree of fairness of each dataset.

### 3.1 Making ENERXICO data Findable

ENERXICO datasets suited for publication will be easily citable and easily findable with the assignation of Persistent Identifiers.

- The codes will be stored in repositories which permit versioning and tags for the identification of official releases and the connection with their outputs.
- Whenever possible, a rich metadata model and the register in disciplinary

repositories will be used to allow other scientists to find the datasets produced by the project.

- Given the variety of the data of the project, the specific solutions and data models adopted for each dataset and software will be found in the corresponding sheet of the DMP.

### 3.2 Making ENERXICO data openly Accessible

The open-data will be made accessible as follows:

- The source-codes of the parallel scientific codes in ENERXICO and related software components licensed as open-source can be included in a web repository for codes and toolkits at the end of the project. It will be useful for archiving project results.
- Datasets access will depend on the different case and it will be described in the corresponding dataset sheet. Restriction of access will be guaranteed in cases ethical or proprietary issues arise. Metadata will be made available as soon and as long as possible, independently on the accessibility of data.

### 3.3 Making ENERXICO data Interoperable

The choice of metadata standards and of the way to access the data is still under discussion between the consortium members. Whenever possible, data coming from other resources will also be described in the DMP. Metadata standards will be chosen to guarantee the maximum interoperability.

### 3.4 Increase ENERXICO data Re-use

The ENERXICO open-datasets will be licensed under Creative Commons data licensing (see Table 5) to let the widest reuse possible of it, since this license allows both commercial and non-commercial use of the data without any restriction. If necessary, an embargo in the data may exist to guarantee publication of results for a maximum of 1 year after the conclusion of the project. In any case, this will be specified in the corresponding dataset sheet.

		Allowed			
Creative Commons	Description	Modification of the content	Commercial use	Free cultural works	Open definition
CC0	Free content, no restrictions	yes	yes	yes	yes
BY	Attribution	yes	yes	yes	yes
BY-SA	Attribution + Share Alike	yes	yes	yes	yes
BY-NC	Non Commercial	yes	no	no	no

BY-ND	No Derivatives	no	yes	no	no
BY-NC-SA		yes	no	no	no
BY-NC-ND		no	no	no	no

Table 5. Data licensing options.

## 4. Allocation of Resource

There is no additional cost for making the ENERXICO datasets identified in Section 2 FAIR:

- The code performance evaluation datasets will be maintained at BSC facilities and included in publications.
- The rest of the open-data will be stored at the project site for at least three years after the end of the project. The infrastructure and personnel funds granted from the European Community will cover the storage, hardware and staff time to manage the servers on which the data will be stored.

## 5. Data security

Each dataset will be evaluated separately and exceptional security measures will be identified and applied. Regular backups for preventing loss of information will be used.

## 6. Ethical aspects

Early warning and hazard assessment can potentially have ethical implications: the diffusion of hazard results or a warning message can be risky for public order, and have social and economic impact, the project will distribute the ultimate simulations results (application to real cases) only under specific conditions and to the appropriate stakeholders, while the scientific results and products, like the models, etc., will be openly accessible. The limitations and conditions of distribution of each dataset will be indicated in the corresponding dataset sheet.

## 7. Annex I

Dataset Sheets

DataSet Sheet (WRF)		Valid values: 1 (totally compliant), 0.5 (partially/ongoing), 0 (not compliant)	
Name	Data simulated with and from the WRF code	Value	F FINDABLE
Description	A set of problems on weather and climate that can be solved where the topography is a key component	0,5	F.1 Persistent identifiers
Data Category	1.01, 1.02, 1.03, and 1.04	0,2	F.2 Rich metadata
Licence	TBD	0	F.3 Data registered in
Repository location	ERAS for input data; To be implemented at ciemates domain for output data	0	F.4 Metadata specify the
Author	Several (Jorge Navarro, Pablo Garcia Müller)		0,175
Naming Conventions	Just structured in a test per folder basis	Value	A ACCESSIBLE
Versioning	No versioning established	0	A.1 Retrievable by the PDI
Format	To be decided, each dataset has its own format	0,5	A.2 Protocol is open, free
Size	Several files (less than 100 files, up to 2 GB per year), under construction	1	A.3 Protocol allows
Storage	TBD	0,3	A.4 Metadata accessible
Archive path	TBD		0,45
Associated metadata	none, only README files	Value	I INTEROPERABLE
Provenance	not structured, different datasets have different origins	0,8	I.1 Language are formal,
Backups needs	None (under Netapp protocol)	0,8	I.2 Vocabulary is FAIR
access permissions	Specific permission given by CEMAT, might be open upon agreement by the authors	0,2	I.3 Metadata includes
Legal/ethical restrictions	None foreseen		0,6
Reproducibility	Possible	Value	R REUSABLE
Data transfer needs	No need, everything is stored together	0,5	R.1 (Meta)data have
Long term preservation	Initially, until the project ends	0,5	R.2 Released with a clear
Metadata management	NA	0,5	R.3 Provenance information
Resources need	NA	0,5	R.4 Domain-relevant
References to other datasets	NA		0,5

DataSet Sheet Topography suite		Valid values: 1 (totally compliant), 0.5 (partially/ongoing), 0 (not compliant)	
Name	Topography Test Suite	Value	F FINDABLE
Description	A set of problems that can be solved where the topography is a key component	0,5	F.1 Persistent identifiers
Data Category	Models and synthesized data (1.01 and 1.03)	0,2	F.2 Rich metadata
Licence	TBD	0	F.3 Data registered in
Repository location	<a href="https://gitlab.lrz.de/sebastian-wolf/topography-benchmarks">https://gitlab.lrz.de/sebastian-wolf/topography-benchmarks</a>	0	F.4 Metadata specify the
Author	Many (Sebastian Wolf, Otilio Rojas, Josep de la Puente, ...)		0,175
Naming Conventions	Just structured in a test per folder basis	Value	A ACCESSIBLE
Versioning	No versioning established	0	A.1 Retrievable by the PDI
Format	To be decided, each dataset has its own format	0,5	A.2 Protocol is open, free
Size	Several files (less than 100 files, 1.1 GB currently), under construction	1	A.3 Protocol allows
Storage	gitlab	0,3	A.4 Metadata accessible
Archive path	gitlab		0,45
Associated metadata	none, only README files	Value	I INTEROPERABLE
Provenance	not structured, different datasets have different origins	0,8	I.1 Language are formal,
Backups needs	None (under git)	0,8	I.2 Vocabulary is FAIR
access permissions	Specific permission given by LRZ, might be open upon agreement by the authors	0,2	I.3 Metadata includes
Legal/ethical restrictions	None foreseen		0,6
Reproducibility	In principle possible	Value	R REUSABLE
Data transfer needs	No need, everything is stored together	0,5	R.1 (Meta)data have
Long term preservation	Initially, until the project ends	0,5	R.2 Released with a clear
Metadata management	NA	0,5	R.3 Provenance information
Resources need	NA	0,5	R.4 Domain-relevant
References to other datasets	NA		0,5

## Software Sheet

Software Sheet - BSIT	
Reference name of the program or	Barcelona Subsurface Imaging Tools (BSIT)
Description	An HPC geophysical imaging tools that includes seismic and EM capabilities for modelling, migration and inversion of 2D or 3D datasets
Author	BSC (several authors inc. Mauricio Hanzich, Josep de la Puente, J.E. Rodriguez, N. Gutierrez, J. Kormann, M. Ferrer, A. Farrés)
Programming language	C
Rules and best coding practices	Internal best practices and documentation rules
Access permissions and license	Intellectual property of Repsol. Pending academic usage license at BSC
Code size	100,000 lines
Repository type	Internal svn repository at BSC, unknown repositories at Repsol
Repository structure	One yearly deliverable in 2010-2018, internally several development branches
Provenance information	Build information on executables
Backup and Archiving needs	None
Legal/ethical restrictions	None
Versioning control and	Not available for public usage
Code transfer needs and security	Private code, under industrial secret
Long term preservation needs	NA
Documentation and inline comments	Internal documentation rules, documentation compiled with doxygen
Metadata management	NA
Resources need	Requirements for sw at each step of the lifecycle (access to repository, computational needs, accessibilities, permissions, ...)

### Software Sheet - WRF

<i>Reference name of the program or</i>	Weather Research Forecast (WRF)
<i>Description</i>	A mesoscale numerical weather prediction system designed for both atmospheric research and operational forecasting applications
<i>Author</i>	NCAR (A couple of modules designed and implemented by CIEMAT staff. PBL module with YSU and SL
<i>Programming language</i>	F90, F95, C
<i>Rules and best coding practices</i>	Internal best practices and documentation rules
<i>Access permissions and license</i>	Open source
<i>Code size</i>	100,000 lines
<i>Repository type</i>	<a href="https://www2.mmm.ucar.edu/wrf/users/downloads.html">https://www2.mmm.ucar.edu/wrf/users/downloads.html</a>
<i>Repository structure</i>	To be implemented at ciemat.es domain
<i>Provenance information</i>	Build information on executables
<i>Backup and Archiving needs</i>	None
<i>Legal/ethical restrictions</i>	None
<i>Versioning control and</i>	wrf4.1
<i>Code transfer needs and security</i>	Open source
<i>Long term preservation needs</i>	NA
<i>Documentation and inline comments</i>	<a href="https://www2.mmm.ucar.edu/wrf/users/pub-doc.html">https://www2.mmm.ucar.edu/wrf/users/pub-doc.html</a>
<i>Metadata management</i>	NA
<i>Resources need</i>	FORTRAN 90 or 95 and C, perl, OpenMP, MPI RSL-LITE, C-shell and Bourne shell, make, M4, sed, awk, netCDF-4, PHD5, Grib-1, NCL, RIB4, ARWboost

### Software Sheet - Alya

<i>Reference name of the program or</i>	Alya
<i>Description</i>	A multi-physics software developed in CASE department designed to work in massive parallel environments.
<i>Author</i>	The main architect is Guillaume Houzeaux, in the context of ENERXICO main contributions are from Matias Avila, Daniel Mira, Ambrus Both and Oriol Lehmkuhl. NCAR (A couple of modules designed and implemented by CIEMAT staff. PBL module with YSU and SL implementations)
<i>Programming language</i>	F90
<i>Rules and best coding practices</i>	Internal best practices and documentation rules
<i>Access permissions and license</i>	Via LUL agreement between BSC and the user
<i>Code size</i>	700,000 lines
<i>Repository type</i>	Git Lab
<i>Repository structure</i>	
<i>Provenance information</i>	Build information on executables
<i>Backup and Archiving needs</i>	None
<i>Legal/ethical restrictions</i>	None
<i>Versioning control and</i>	Following git format
<i>Code transfer needs and security</i>	Open source
<i>Long term preservation needs</i>	NA
<i>Documentation and inline comments</i>	<a href="https://gitlab.bsc.es/alya/alya/-/wikis/home">https://gitlab.bsc.es/alya/alya/-/wikis/home</a>
<i>Metadata management</i>	NA
<i>Resources need</i>	FORTRAN 90 or 95 and C, MPI, make, python, paraview

### Software Sheet - SeisSol

<i>Reference name of the program or</i>	SeisSol
<i>Description</i>	SeisSol solves seismic wave propagation (elastic, viscoelastic) and dynamic rupture problems on heterogeneous 3D models.
<i>Author</i>	LMU and TUM
<i>Programming language</i>	C++, Fortran, Python
<i>Rules and best coding practices</i>	Internal best practices and documentation rules
<i>Access permissions and license</i>	Open source (BSD-3)
<i>Code size</i>	950000
<i>Repository type</i>	<a href="https://github.com/SeisSol/SeisSol">https://github.com/SeisSol/SeisSol</a>
<i>Repository structure</i>	Rolling release, snapshots for verification, development/feature branches
<i>Provenance information</i>	Build information on executables
<i>Backup and Archiving needs</i>	None
<i>Legal/ethical restrictions</i>	None
<i>Versioning control and</i>	feature branches with pull requests/code review
<i>Code transfer needs and security</i>	Open source
<i>Long term preservation needs</i>	NA
<i>Documentation and inline comments</i>	<a href="https://seissol.readthedocs.io/en/latest/">https://seissol.readthedocs.io/en/latest/</a>
<i>Metadata management</i>	NA
<i>Resources need</i>	gcc or intel, numpy, pamefis, libxsmm, pspamm, mpi, netcdf, hdf5, external meshing tool (gmsh, simmetrix)

### Software Sheet - ExaHyPE

<i>Reference name of the program or workflow</i>	ExaHyPE Engine
<i>Description</i>	An engine for the simulation of systems of hyperbolic PDEs, as stemming from conservation laws
<i>Author</i>	The ExaHyPE Consortium
<i>Programming language</i>	C++, Fortran, Python
<i>Rules and best coding practices</i>	Internal best practices and documentation rules
<i>Access permissions and license</i>	Open source
<i>Code size</i>	500,000 lines of code
<i>Repository type</i>	<a href="https://gitlab.lrz.de/exahype/ExaHyPE-Engine">https://gitlab.lrz.de/exahype/ExaHyPE-Engine</a>
<i>Repository structure</i>	Running release mode, stable versions available on branches, several development branches
<i>Provenance information</i>	Build information on executables
<i>Backup and Archiving needs</i>	None
<i>Legal/ethical restrictions</i>	None
<i>Versioning control and rules/workflows managing</i>	Running release mode
<i>Code transfer needs and security</i>	Open source
<i>Long term preservation needs</i>	NA
<i>Documentation and inline comments rules</i>	<a href="http://www.peano-framework.org/exahype/guidebook.pdf">http://www.peano-framework.org/exahype/guidebook.pdf</a>
<i>Metadata management</i>	NA
<i>Resources need</i>	Intel or GNU compilers, MPI, TBB, Python3

### Software Sheet - SEM46

<i>Reference name of the program or workflow</i>	Barcelona Subsurface Imaging Tools (BSIT)
<i>Description</i>	3D seismic modeling and inversion code, developed mainly in the frame of the SEISCOPE project ( <a href="https://seiscope2.osug.fr/">https://seiscope2.osug.fr/</a> ) for tackling modeling and full waveform inversion topics from the near surface to the deep crustal scale.
<i>Author</i>	Romain Brossier <a href="mailto:romain.brossier@univ-grenoble-alpes.fr">romain.brossier@univ-grenoble-alpes.fr</a>
<i>Programming language</i>	Internal best practices and documentation rules
<i>Rules and best coding practices</i>	Internal best practices and documentation rules
<i>Access permissions and license</i>	SEM46 is available upon request to Romain Brossier <a href="mailto:romain.brossier@univ-grenoble-alpes.fr">romain.brossier@univ-grenoble-alpes.fr</a> SEM46 relies on a BSD-like license but include restrictions for the diffusion, imposed by the funding partners of the SEISCOPE project
<i>Code size</i>	
<i>Repository type</i>	
<i>Repository structure</i>	
<i>Provenance information</i>	
<i>Backup and Archiving needs</i>	None
<i>Legal/ethical restrictions</i>	None
<i>Versioning control and rules/workflows managing</i>	
<i>Code transfer needs and security</i>	
<i>Long term preservation needs</i>	NA
<i>Documentation and inline comments rules</i>	
<i>Metadata management</i>	NA
<i>Resources need</i>	

### Software Sheet - DualSPHysics and Black Hole codes

<i>Reference name of the program or workflow</i>	DualSPHysics and Black Hole codes
<i>Description</i>	DualSPHysics is a hardware accelerated Smoothed Particle Hydrodynamics code developed to solve free-surface flow problems. The extension to multiphase flow has been developed through the ENERXICO project and is called the Black Hole code.
<i>Author</i>	The DualSPHysics and Black Hole codes developers
<i>Programming language</i>	C++, CUDA, Java
<i>Rules and best coding practices</i>	Internal best practices and documentation rules
<i>Access permissions and license</i>	Open source (DualSPHysics), non-open source (BH code)
<i>Code size</i>	100,000 lines of code
<i>Repository type</i>	<a href="https://github.com/DualSPHysics/DualSPHysics/tree/develop">https://github.com/DualSPHysics/DualSPHysics/tree/develop</a>
<i>Repository structure</i>	Running release mode, stable versions available on branches, several development branches
<i>Provenance information</i>	Build information on executables
<i>Backup and Archiving needs</i>	None
<i>Legal/ethical restrictions</i>	None
<i>Versioning control and rules/workflows managing</i>	Running release mode
<i>Code transfer needs and security</i>	Open source
<i>Long term preservation needs</i>	NA
<i>Documentation and inline comments rules</i>	<a href="https://dual.sphysics.org/doxygen/">https://dual.sphysics.org/doxygen/</a>
<i>Metadata management</i>	NA
<i>Resources need</i>	C++ and CUDA compilers, MPI, Java

Formats definition

## Data Categories definition

Data Categories Definition							
Find your dataset category from the listed below. If you need a new sub-category, you can add more at the end of the corresponding category list.							
<b>1 Scientific Data</b>		<b>2 Software</b>		<b>3 Administrative docs</b>		<b>4 Other</b>	
Name	Definition	Name	Definition	Name	Definition	Name	Definition
<b>1.01 metadata</b>		<b>2.01 libraries</b>		<b>3.01 Documents</b>	Any documentation, either	<b>4.01 metadata</b>	Any data describing data
<b>1.02 experimental data</b>	Data coming from	<b>2.02 applications</b>		<b>3.02 Internal reports</b>	Meeting minutes, Internal		
<b>1.03 synthetic data</b>	Data produced from	<b>2.03 services</b>		<b>3.03 Deliverables</b>	Project outputs documents		
<b>1.04 test data</b>	Data used to test software	<b>2.04 APIs - source code</b>					

## Data repositories Information

Data Repositories Information					
Data repositories are listed below. If you need, you can add more at the end of the list.					
Repository	Extended name	Location	URL	Permissions	...

## DataSet Sheet Template - do not modify

DataSet Sheet Template - do not modify!		Valid values: 0 (totally compliant), 0.5 (partially/coming), 0 (not)
Name	Descriptive name to identify the dataset	Value <b>F</b> <b>FINDABLE</b>
Description	Short description of the contents	F.1 Persistent identifiers (PDI)
Data Category	Data category code (see Table Data Category for the corresponding codes)	F.2 Rich metadata
Licence	Chosen among the most appropriated and most open ones	F.3 Data rights listed in searchable metadata
Repository location	Institutional or public repository name and URL, if available	F.4 Metadata specify the PDI
Author	Data author(s) name(s)	
Naming Conventions	File names, structure and conventions	Value <b>A</b> <b>ACCESSIBLE</b>
Versioning	How and where the version of the dataset can be found	A.1 Retrievable by the PDI with a standard
Format	document, General or specific format - libraries or parsing code	A.2 Protocol is open, free
Size	Total or single file size * n. of files	A.3 Protocol allows authentication and
Storage	Physical support	A.4 Metadata accessible beyond the data
Archive path	Folders structure	
Associated metadata	reference to metadata standards	Value <b>I</b> <b>INTEROPERABLE</b>
Provenance	Structured dataset origin information	I.1 Language and format, accessible, shared and
Backup needs	Periodicity, subsets, backup needs analysis, etc.	I.2 Vocabulary is FAIR
Access permissions	Lifecycle dependency: only specific groups of collaborators, all partners, whole community, ...	I.3 Metadata includes qualified references to
Legal/ethical restrictions	Privacy and security issues	
Reproducibility	If yes: connection to code and environment	Value <b>R</b> <b>REUSABLE</b>
Data transfer needs	Replicas and periodic transfers to/from other repositories	R.1 (Metadata) have plurality of accurate and
Long term preservation	Needs at 3-5-7-10 years (if any)	R.2 Released with a clear and accessible data
Metadata management	Way to access metadata when data are not available	R.3 Provenance information
Resources need	Analysis of resources needs at each step of data lifecycle	R.4 Domain-relevant community standards
References to other datasets	If applicable, explain which and why	

## Software Sheet - do not modify

Software Sheet - do not modify	
<i>Reference name of the program or Description</i>	Name of the code Brief description of the functionality and applicability of the software
<i>Author</i>	specify
<i>Programming language</i>	specify the programming language(s)
<i>Rules and best coding practices</i>	conventions for filenames, link to an external manual, if exists (ex PEP8, etc.)
<i>Access permissions and license</i>	lifecycle dependency: only specific groups of collaborators, all partners, whole community, etc.
<i>Code size</i>	if relevant (to be updated)
<i>Repository type</i>	GitHub, GitLab, Bitbucket, SourceForge
<i>Repository structure</i>	Branches, tags, etc.
<i>Provenance information</i>	Containers, virtual environments
<i>Backup and Archiving needs</i>	if any
<i>Legal/ethical restrictions</i>	if any
<i>Versioning control and Code transfer needs and security</i>	Specify the repository if any
<i>Long term preservation needs</i>	Only if applies to a given official release version
<i>Documentation and inline comments</i>	specify
<i>Metadata management</i>	(available even when the software is not)
<i>Resources need</i>	Requirements for sw at each step of the lifecycle (access to repository, computational needs, accessibilities, permissions, ...)

## FAIR DATA definitions - do not modify

FAIR DATA definitions - do not modify					
<a href="#">GoFAIR - principles definition link</a>					
To evaluate the FAIRness of a dataset, a punctuation should be given to each of the following points (elaborated by www.force11.org/valid values: 1 (totally compliant), 0.5 (partially compliant), 0 (not compliant))					
<b>F</b>	<b>FINDABLE</b>	<b>A</b>	<b>ACCESSIBLE</b>	<b>I</b>	<b>INTEROPERABLE</b>
F.1	(Meta)data are assigned a globally unique and persistent identifier	A.1	(Meta)data are retrievable by their identifier using a standardized communications protocol	I.1	(Meta)data use a formal, machine-actionable vocabulary (e.g. RDF)
F.2	Data are described with rich metadata	A.2	The protocol is open, free, and universal	I.2	Metadata use vocabularies that are interoperable with other commonly used vocabularies
F.3	(Meta)data are registered or indexed	A.3	The protocol allows for authentication and authorization	I.3	(Meta)data include qualified statements of provenance
F.4	Metadata clearly and explicitly include the identifier of the data to which they apply	A.4	Metadata are accessible, even when the primary data are not		
				<b>R</b>	<b>REUSABLE</b>
				R.1	(Meta)data are richly described with machine-actionable information
				R.2	(Meta)data are released with an open license
				R.3	(Meta)data are associated with data or services
				R.4	(Meta)data meet domain-specific requirements

## Data licensing

Creative Commons	Description	Modification of the work	Allowed		
			Commercial Use	Free cultural works	Open definition
CC0	Free content, no restrictions	yes	yes	yes	yes
BY	Attribution	yes	yes	yes	yes
BY-SA	Attribution+ShareAlike	yes	yes	yes	yes
BY-NC	NonCommercial	yes	no	no	no
BY-ND	NoDerivatives	no	yes	no	no
BY-NC-SA		yes	no	no	no
BY-NC-ND		no	no	no	no

### Definitions of common license types

data.world: Common license types for datasets

[For details, click here](#)

<b>Public Domain</b>	The work has been dedicated to the public domain by waiving all rights to the work worldwide under copyright law, including all related and neighboring rights, to the extent allowed by law. You must give appropriate credit, provide a link to the license, and indicate if changes were made. You may do so in any reasonable manner, but not in any way that suggests the licensor endorses you or your use.
<b>Attribution</b>	if you remix, transform, or build upon the material, you must distribute your contributions under the same license as the original.
<b>Share-alike</b>	
<b>Non-commercial</b>	You may not use the material for commercial purposes.
<b>Database Only</b>	License applies to the database only and not its contents or data.
<b>No Derivatives</b>	No Derivative Works. You may not alter, transform, or build upon this work.