

EGU21-12384, updated on 16 Sep 2021  
<https://doi.org/10.5194/egusphere-egu21-12384>  
EGU General Assembly 2021  
© Author(s) 2021. This work is distributed under  
the Creative Commons Attribution 4.0 License.



## AI-Ready Training Datasets for Earth Observation: Enabling FAIR data principles for EO training data.

**Alastair McKinstry**<sup>1,2</sup>, Oisin Boydell<sup>3,4</sup>, Quan Le<sup>3,4</sup>, Inder Preet<sup>3,4</sup>, Jennifer Hanafin<sup>1,2</sup>, Manuel Fernandez<sup>1,2</sup>, Adam Warde<sup>1,2</sup>, Venkatesh Kannan<sup>1,2</sup>, and Patrick Griffiths<sup>5</sup>

<sup>1</sup>NUI Galway, Irish Centre for High-End Computing, Galway, Ireland

<sup>2</sup>Irish Centre for High End Computing, NUI Galway, Galway, Ireland

<sup>3</sup>University College Dublin, Belfield, Dublin 4, Ireland

<sup>4</sup>CeADAR, UCD, Ireland

<sup>5</sup>ESA ESRIN, Via Galileo Galilei, 1, 00044 Frascati RM, Italy

The ESA-funded AIREO project [1] sets out to produce AI-ready training dataset specifications and best practices to support the training and development of machine learning models on Earth Observation (EO) data. While the quality and quantity of EO data has increased drastically over the past decades, availability of training data for machine learning applications is considered a major bottleneck. The goal is to move towards implementing FAIR data principles for training data in EO, enhancing especially the finability, interoperability and reusability aspects. To achieve this goal, AIREO sets out to provide a training data specification and to develop best practices for the use of training datasets in EO. An additional goal is to make training data sets self-explanatory (“AI-ready”) in order to expose challenging problems to a wider audience that does not have expert geospatial knowledge.

Key elements that are addressed in the AIREO specification are granular and interoperable metadata (based on STAC), innovative Quality Assurance metrics, data provenance and processing history as well as integrated feature engineering recipes that optimize platform independence. Several initial pilot datasets are being developed following the AIREO data specifications. These pilot applications include for example forest biomass, sea ice detection and the estimation of atmospheric parameters. An API for the easy exploitation of these datasets will be provided to allow the Training Datasets (TDS) to work against EO catalogs (based on OGC STAC catalogs and best practises from ML community) to allow updating and updated model training over time.

This presentation will present the first version of the AIREO training dataset specification and will showcase some elements of the best-practices that were developed. The AIREO compliant pilot datasets will be presented which are openly accessible and community feedback is explicitly encouraged.

[1] <https://aireo.net/>

