

Open Geospatial Consortium

Submission Date: <2023-03-07>

Approval Date: <yyyy-mm-dd>

Publication Date: <yyyy-mm-dd>

External identifier of this OGC® document: <http://www.opengis.net/doc/{doc-type}/{standard}{m.n}>

Internal reference number of this OGC® document: 23-008r1

Version: 1.0.0

Editor: Peng Yue, Boyi Shangguan

OGC Training Data Markup Language for Artificial Intelligence (TrainingDML-AI) Part1: Conceptual Model Standard

Copyright notice

Copyright © <year> Open Geospatial Consortium

To obtain additional rights of use, visit <http://www.ogc.org/legal/>

Warning

This document is not an OGC Standard. This document is distributed for review and comment. This document is subject to change without notice and may not be referred to as an OGC Standard.

Recipients of this document are invited to submit, with their comments, notification of any relevant patent rights of which they are aware and to provide supporting documentation.

Document type: OGC® Standard

Document subtype: Conceptual Model

Document stage: Draft

Document language: English

License Agreement

Permission is hereby granted by the Open Geospatial Consortium, ("Licensor"), free of charge and subject to the terms set forth below, to any person obtaining a copy of this Intellectual Property and any associated documentation, to deal in the Intellectual Property without restriction (except as set forth below), including without limitation the rights to implement, use, copy, modify, merge, publish, distribute, and/or sublicense copies of the Intellectual Property, and to permit persons to whom the Intellectual Property is furnished to do so, provided that all copyright notices on the intellectual property are retained intact and that each person to whom the Intellectual Property is furnished agrees to the terms of this Agreement.

If you modify the Intellectual Property, all copies of the modified Intellectual Property must include, in addition to the above copyright notice, a notice that the Intellectual Property includes modifications that have not been approved or adopted by LICENSOR.

THIS LICENSE IS A COPYRIGHT LICENSE ONLY, AND DOES NOT CONVEY ANY RIGHTS UNDER ANY PATENTS THAT MAY BE IN FORCE ANYWHERE IN THE WORLD.

THE INTELLECTUAL PROPERTY IS PROVIDED "AS IS", WITHOUT WARRANTY OF ANY KIND, EXPRESS OR IMPLIED, INCLUDING BUT NOT LIMITED TO THE WARRANTIES OF MERCHANTABILITY, FITNESS FOR A PARTICULAR PURPOSE, AND NONINFRINGEMENT OF THIRD PARTY RIGHTS. THE COPYRIGHT HOLDER OR HOLDERS INCLUDED IN THIS NOTICE DO NOT WARRANT THAT THE FUNCTIONS CONTAINED IN THE INTELLECTUAL PROPERTY WILL MEET YOUR REQUIREMENTS OR THAT THE OPERATION OF THE INTELLECTUAL PROPERTY WILL BE UNINTERRUPTED OR ERROR FREE. ANY USE OF THE INTELLECTUAL PROPERTY SHALL BE MADE ENTIRELY AT THE USER'S OWN RISK. IN NO EVENT SHALL THE COPYRIGHT HOLDER OR ANY CONTRIBUTOR OF INTELLECTUAL PROPERTY RIGHTS TO THE INTELLECTUAL PROPERTY BE LIABLE FOR ANY CLAIM, OR ANY DIRECT, SPECIAL, INDIRECT OR CONSEQUENTIAL DAMAGES, OR ANY DAMAGES WHATSOEVER RESULTING FROM ANY ALLEGED INFRINGEMENT OR ANY LOSS OF USE, DATA OR PROFITS, WHETHER IN AN ACTION OF CONTRACT, NEGLIGENCE OR UNDER ANY OTHER LEGAL THEORY, ARISING OUT OF OR IN CONNECTION WITH THE IMPLEMENTATION, USE, COMMERCIALIZATION OR PERFORMANCE OF THIS INTELLECTUAL PROPERTY.

This license is effective until terminated. You may terminate it at any time by destroying the Intellectual Property together with all copies in any form. The license will also terminate if you fail to comply with any term or condition of this Agreement. Except as provided in the following sentence, no such termination of this license shall require the termination of any third party end-user sublicense to the Intellectual Property which is in force as of the date of notice of such termination. In addition, should the Intellectual Property, or the operation of the Intellectual Property, infringe, or in LICENSOR's sole opinion be likely to infringe, any patent, copyright, trademark or other right of a third party, you agree that LICENSOR, in its sole discretion, may terminate this license without any compensation or liability to you, your licensees or any other party. You agree upon termination of any kind to destroy or cause to be destroyed the Intellectual Property together with all copies in any form, whether held by you or by any third party.

Except as contained in this notice, the name of LICENSOR or of any other holder of a copyright in all or part of the Intellectual Property shall not be used in advertising or otherwise to promote the sale, use or other dealings in this Intellectual Property without prior written authorization of LICENSOR or such copyright holder. LICENSOR is and shall at all times be the sole entity that may authorize you or any third party to use certification marks, trademarks or other special designations to indicate compliance with any LICENSOR standards or specifications. This Agreement is governed by the laws of the Commonwealth of Massachusetts. The application to this Agreement of the United Nations Convention on Contracts for the International Sale of Goods is hereby expressly excluded. In the event any provision of this Agreement shall be deemed unenforceable, void or invalid, such provision shall be modified so as to make it valid and enforceable, and as so modified the entire Agreement shall remain in full force and effect. No decision, action or inaction by LICENSOR shall be construed to be a waiver of any rights or remedies available to it.

Table of Contents

1. Scope	8
2. Conformance	9
3. Normative References	10
4. Terms and Definitions	11
4.1. Artificial Intelligence (AI)	11
4.2. Machine Learning (ML)	11
4.3. Deep Learning (DL)	11
4.4. Training Dataset	11
4.5. Label	11
4.6. Provenance	12
4.7. Quality	12
4.8. Earth Observation	12
4.9. Scene Classification	12
4.10. Object Detection	12
4.11. Semantic Segmentation	12
4.12. Change Detection	12
4.13. 3D Model Reconstruction	12
4.14. Generative model	13
5. Conventions	14
5.1. Identifiers	14
5.2. Abbreviated terms	14
5.3. UML Notation	14
6. Overview	16
6.1. AI tasks for EO	16
6.2. Modularization	17
6.3. General modeling principles	17
6.3.1. Element modeling	17
6.3.2. Class Hierarchy and Inheritance of Properties and Relations	18
6.3.3. Definition of the Semantics for all Classes, Properties, and Relations	18
6.3.4. Data Integrity, Authenticity, and Non-repudiation	18
6.4. Extending TrainingDML-AI	18
7. TrainingDML-AI UML Model	20
7.1. ISO dependencies	20
7.2. Overview of the UML model	21
7.3. AI_TrainingDataset	22
7.3.1. Requirements	23
7.3.2. Class definitions	24
7.4. AI_TrainingData	25

7.4.1. Requirements	26
7.4.2. Class definitions	27
7.5. AI_Task	27
7.5.1. Requirements	28
7.5.2. Class definitions	29
7.6. AI_Label	29
7.6.1. Requirements	30
7.6.2. Class definitions	31
7.7. AI_Labeling	31
7.7.1. Requirements	32
7.7.2. Class definitions	33
7.8. AI_TDChangeset	33
7.8.1. Requirements	34
7.8.2. Class definitions	35
7.9. AI_DataQuality	35
7.9.1. Requirements	36
7.9.2. Class definitions	37
8. TrainingDML-AI Data Dictionary	38
8.1. ISO Classes	38
8.1.1. Feature (from ISO 19107:2019)	38
8.1.2. MD_Band (from ISO 19115-1:2014)	38
8.1.3. MD_Scope (from ISO 19115-1:2014)	38
8.1.4. EX_Extent (from ISO 19115-1:2014)	38
8.1.5. CI_Citation (from ISO 19115-1:2014)	39
8.1.6. DataQuality (from ISO 19157-1)	39
8.1.7. QualityElement (from ISO 19157-1)	39
8.2. AI_TrainingDataset	39
8.2.1. Classes	39
8.3. AI_TrainingData	41
8.3.1. Classes	41
8.4. AI_Task	43
8.4.1. Classes	43
8.5. AI_Label	43
8.5.1. Classes	44
8.6. AI_Labeling	45
8.6.1. Classes	45
8.7. AI_TDChangeset	46
8.7.1. Classes	47
8.8. AI_DataQuality	47
8.8.1. Classes	47
Appendix A: Abstract Test Suite (Normative)	49

A.1. Introduction	49
A.2. Conformance class AI_TrainingDataset	49
A.3. Conformance class AI_TrainingData	50
A.4. Conformance class AI_Task	51
A.5. Conformance class AI_Label	52
A.6. Conformance class AI_Labeling	53
A.7. Conformance class AI_TDChangeset	54
A.8. Conformance class AI_DataQuality	55
Appendix B: Example (Informative)	57
B.1. TrainingDataset encoding examples	57
B.1.1. WHU-RS19 dataset	57
B.1.2. DOTA-v1.5 dataset	57
B.1.3. KITTI 2D object detection dataset	57
B.1.4. GID dataset	58
B.1.5. Toronto3D dataset	58
B.1.6. WHU-Building dataset	58
B.1.7. California change detection dataset	58
B.1.8. WHU MVS dataset	59
B.2. DataQuality encoding example	59
B.2.1. WHU-RS19 data quality	59
B.3. TDChangeset encoding examples	59
B.3.1. DOTA-v1.5 changeset	59
Appendix C: Revision History (Informative)	60
Appendix D: Bibliography	61

i. Abstract

The Training Data Markup Language for Artificial Intelligence (TrainingDML-AI) Standard aims to develop the UML model and encodings for geospatial machine learning training data. Training data plays a fundamental role in Earth Observation (EO) Artificial Intelligence Machine Learning (AI/ML), especially Deep Learning (DL). It is used to train, validate, and test AI/ML models. This Standard defines a UML model and encodings consistent with the OGC Standards baseline to exchange and retrieve the training data in the Web environment.

The TrainingDML-AI Standard provides detailed metadata for formalizing the information model of training data. This includes but is not limited to the following aspects:

- How the training data is prepared, such as provenance or quality;
- How to specify different metadata used for different ML tasks such as scene/object/pixel levels;
- How to differentiate the high-level training data information model and extended information models specific to various ML applications;
- How to introduce external classification schemes and flexible means for representing ground truth labeling.

ii. Keywords

The following are keywords to be used by search engines and document catalogues.

ogcdoc, OGC document, artificial intelligence, machine learning, deep learning, earth observation, remote sensing, training data, training sample, UML

iii. Preface

Attention is drawn to the possibility that some of the elements of this document may be the subject of patent rights. The Open Geospatial Consortium shall not be held responsible for identifying any or all such patent rights.

iv. Security Considerations

No security considerations have been made for this Standard.

v. Submitting organizations

The following organizations submitted this Document to the Open Geospatial Consortium (OGC):

- Wuhan University
- Pixalytics Ltd
- National Geospatial-Intelligence Agency
- George Mason University
- Laboratoire d’Informatique de Grenoble
- South Digital Technology Co., Ltd

- Wuhan University of Technology
- Hubei University
- Chongqing Changan Automobile Co., Ltd

vi. Submitters

All questions regarding this submission should be directed to the editor or the submitters:

Name	Affiliation
Peng Yue	Wuhan University
Jianya Gong	Wuhan University
Ruixiang Liu	Wuhan University
Dayu Yu	Wuhan University
Samantha Lavender	Pixalytics Ltd
Jim Antonisse	National Geospatial-Intelligence Agency
Liping Di	George Mason University
Eugene Yu	George Mason University
Danielle Ziébelin	Laboratoire d'Informatique de Grenoble
Boyi Shangguan	South Digital Technology Co., Ltd
Lei Hu	South Digital Technology Co., Ltd
Liangcun Jiang	Wuhan University of Technology
Mingda Zhang	Hubei University
Kai Yan	Chongqing Changan Automobile Co., Ltd

vii. Acknowledgments

Thanks to the members of the TrainingDML-AI Standards Working Group of the OGC as well as all contributors of change requests and comments. In particular: Scott Simmons, Carl Reed, Ivana Ivánová, Emily Daemen, Jon Duckworth, Zheheng Liang, Jibo Xie, Yuqi Bai, Winnie Shiu, Ignacio Correas, Chenxiao Zhang, Zhipeng Cao, Haofeng Tan, Yinyin Pan, Hanwen Xu, Shuaiqi Liu, Hao Li, Ming Wang, Kaixuan Wang, Haipeng Deng, Gobe Hobona, Chris Little, Kathi Schleidt, Rodolfo Orozco.

Chapter 1. Scope

Training data is the building block of machine learning models. These models now constitute the majority of machine learning applications in Earth science. Training data is used to train AI/ML models, and to then validate model results. Formalizing and documenting the training data by characterizing the training data content, metadata, data quality, and provenance, and so forth is essential.

This OGC Training Data Standard describes work actions around training data:

- Documents the UML model with a target of maximizing the interoperability and usability of EO imagery training data;
- Defines different AI/ML tasks and labels in earth observation in terms of supervised learning, including scene level, object level and pixel level tasks;
- Describes the description of the permanent identifier, version, license, training data size, measurement or imagery used for annotation, and so on;
- Defines the description of quality (e.g., training data errors, training data unrepresentativeness) and the provenance (labeling, labeler, and labeling procedure).

Chapter 2. Conformance

This TrainingDML-AI Standard defines a conceptual model that is independent of any encoding or formatting technologies. The standardization targets for this Standard is:

- TrainingDML-AI Conceptual Model

Conformance with this Standard shall be checked using all the relevant tests specified in Annex A (normative) of this document. The framework, concepts, and methodology for testing, and the criteria to be achieved to claim conformance are specified in the OGC Compliance Testing Policies and Procedures and [the OGC Compliance Testing web site](#).

All requirements-classes and conformance-classes described in this document are owned by the standard identified.

Chapter 3. Normative References

The following normative documents contain provisions that, through reference in this text, constitute provisions of this document. For dated references, subsequent amendments to, or revisions of, any of these publications do not apply. For undated references, the latest edition of the normative document referred to applies.

- ISO 19103:2015 Geographic information — Conceptual schema language
- ISO 19107:2019 Geographic information — Spatial schema
- ISO 19115-1:2014 Geographic information — Metadata — Part 1: Fundamentals
- ISO 19157-1 Geographic information — Data quality — Part 1: General requirements

Chapter 4. Terms and Definitions

This document used the terms defined in OGC Policy Directive 49, which is based on the ISO/IEC Directives, Part 2, Rules for the structure and drafting of International Standards. In particular, the word “shall” (not “must”) is the verb form used to indicate a requirement to be strictly followed to conform to this Standard and OGC documents do not use the equivalent phrases in the ISO/IEC Directives, Part 2.

For the purposes of this document, the following additional terms and definitions apply.

4.1. Artificial Intelligence (AI)

refers to a set of methods and technologies that can empower machines or software to learn and perform tasks like humans.

4.2. Machine Learning (ML)

is an important branch of artificial intelligence that gives computers the ability to improve their performance without explicitly being programmed to do so. ML processes create models from training data by using a set of learning algorithms, and then can use these models to make predictions. Depending on whether the training data include labels, the learning algorithms can be divided into supervised and unsupervised learning.

4.3. Deep Learning (DL)

is a subset of machine learning, which is essentially a neural network with three or more layers. The number of layers is referred to as depth. While a neural network with a single layer can still make approximate predictions, additional hidden layers can help to optimize and refine for accuracy.

SOURCE: <https://www.ibm.com/topics/deep-learning>

4.4. Training Dataset

a collection of samples, often labelled in terms of supervised learning. A training dataset can be divided into training, validation, and test sets. Training samples are different from samples in OGC Observations & Measurements (O&M). They are often collected in purposive ways that deviate from purely probability sampling, with known or expected results labelled as values of a dependent variable for generating a trained predictive model.

4.5. Label

refers to known or expected results annotated as values of a dependent variable in training samples. A training sample label is different from those on a geographical map, which are known as map labels or annotations.

4.6. Provenance

information about entities, activities, and people involved in producing a piece of data or thing, which can be used to form assessments about its quality, reliability or trustworthiness. In this standard provenance is a record of how training data were prepared.

SOURCE: W3C (<https://www.w3.org/TR/prov-overview/>)

4.7. Quality

degree to which a set of inherent characteristics fulfils requirements [ISO 9000:2005, definition3.1.1]. Quality of training data (such as data imbalance and mislabeling) can impact the performance of AI/ML models.

4.8. Earth Observation

data and information collected about our planet, whether atmospheric, oceanic or terrestrial. This includes space-based or remotely-sensed data, as well as ground-based or in situ data.

SOURCE: GEO (https://earthobservations.org/geo_wwd.php)

4.9. Scene Classification

task of identifying scene categories of images, on the basis of a training set of images whose scene categories are known.

4.10. Object Detection

task of recognizing objects such as cars from images. The objects are often localized using bounding boxes.

4.11. Semantic Segmentation

task of assigning class labels to pixels of images or points of point clouds.

4.12. Change Detection

task that finds the changes in an area between images taken at different times.

4.13. 3D Model Reconstruction

task that builds 3D objects and scenes from multi-view images.

4.14. Generative model

is one of the methods of large model training, which improve model performance through unsupervised pre-training. In the fine-tuning phase, labelled data plays a critical role in optimizing the model for specific vertical domains or tasks. By incorporating labelled data, the model can learn to accurately identify and extract relevant features, leading to better performance on specific downstream tasks. Overall, the combination of generative models and fine-tuning with labelled data can significantly improve the performance of large models in specialized domains or tasks.

Chapter 5. Conventions

This section provides details and examples for any conventions used in the document. Examples of conventions are symbols, abbreviations, use of XML schema, or special notes regarding how to read the document.

5.1. Identifiers

The normative provisions in this specification are denoted by the URI:

<http://www.opengis.net/spec/TrainingDML-AI-1/1.0>

All requirements and conformance tests that appear in this document are denoted by partial URIs which are relative to this base.

5.2. Abbreviated terms

In this document the following abbreviations and acronyms are used or introduced:

- AI — Artificial Intelligence
- DL — Deep Learning
- EO — Earth Observation
- ISO — International Organization for Standardization
- JSON — JavaScript Object Notation
- LC — Land Cover
- LU — Land Use
- ML — Machine Learning
- OGC — Open Geospatial Consortium
- RS — Remote Sensing
- TD — Training Data
- UML — Unified Modelling Language
- XML — Extensible Markup Language

5.3. UML Notation

The Standard is presented in this document through diagrams using the Unified Modeling Language (UML) static structure diagram. The UML notations used in this Standard are described in the diagram in Figure 1.

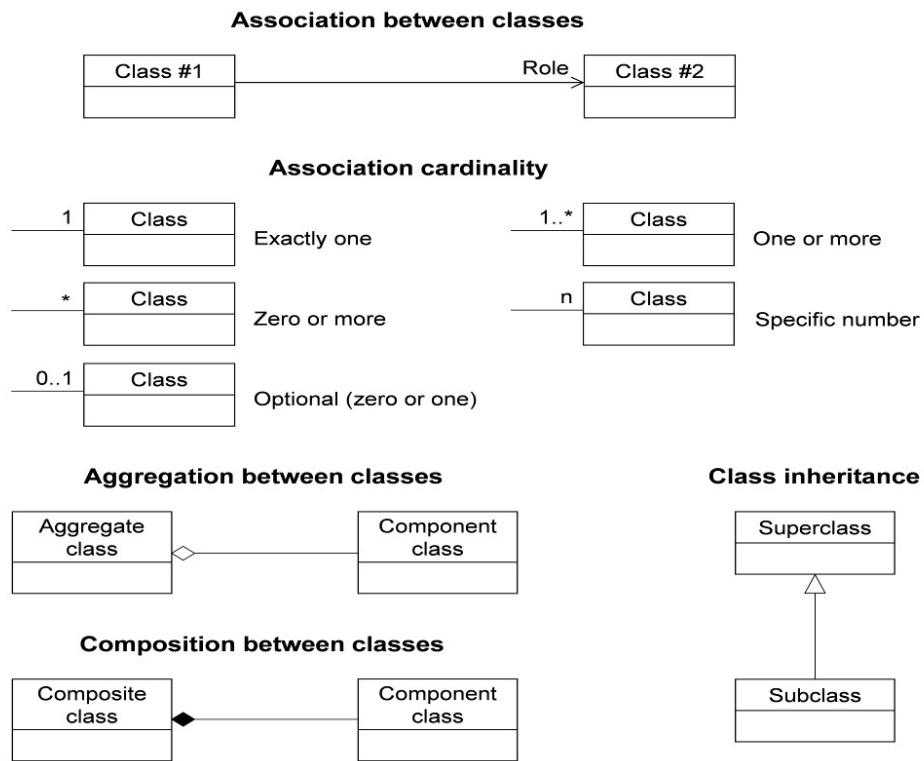


Figure 1. UML notation (see ISO TS 19103, Geographic information—Conceptual schema language).

All associations between model elements in the TrainingDML-AI Conceptual Model are unidirectional. Thus, associations in the model are navigable in only one direction. The direction of navigation is depicted by an arrowhead. In general, the context an element takes within the association is indicated by its role. The role is displayed near the target of the association. If the graphical representation is ambiguous though, the position of the role has to be drawn to the element the association points to.

The following stereotypes are used in this model.

- «DataType» defines a set of properties that lack identity. A data type is a classifier with no operations, whose primary purpose is to hold information.
- «CodeList» enumerates the valid attribute values. In contrast to Enumeration, the list of values is open and, thus, not given inline in the TrainingDML-AI UML Model. The allowed values can be provided within an external code list.

Chapter 6. Overview

The TrainingDML-AI Conceptual Model Standard defines how to represent and exchange ML training data. The conceptual model includes the most relevant training data entities from datasets, to instances (i.e. individual training samples), to labels. The conceptual schema specifies how and into which parts of the training data should be decomposed and classified.

The TrainingDML-AI conceptual model (Clause 7) is formally specified using UML class diagrams, complemented by a data dictionary (Clause 8) providing the definitions and explanations of the object classes and attributes. This conceptual model provides the basis for specifying encoding implemented in languages such as JSON, or XML.

6.1. AI tasks for EO

In recent years AI/ML is increasingly used in the EO domain. The new AI/ML algorithms frequently require large training datasets as benchmarks. AI/ML TD have been used in many EO applications to calibrate the performance of AI/ML models. Many efforts have been made to produce training datasets to make accurate predictions. As a result, a number of training datasets are publicly available, with new datasets being constantly released. In the EO domain, examples of AI/ML training datasets have been developed in various tasks including the following typical scenarios:

- Scene classification. These algorithms determine image categories from numerous pictures (e.g., agricultural, forest, and beach scenes). The training samples are a series of labelled pictures. The data can be either from satellite, drones, or aircrafts. The metadata of the datasets often includes the number of training samples, the number of classes, and the image size.
- Object detection. These algorithms detect and localize different objects (e.g., airplanes, cars and building) in a single image. The image can be optical or non-optical, such as Synthetic Aperture Radar (SAR). Recent work also suggests an increasing focus on object detection from street view imagery. Objects can be labelled with two forms of bounding boxes, i.e., oriented and horizontal bounding boxes. The geometry of a bounding box can be expressed using top-left/bottom-right coordinates, coordinates of four corners, or center coordinates along with the length and width of the box.
- Semantic segmentation. In terms of Land cover (LC) and land use (LU) classification, this process assigns a LC/LU class label to a pixel (or groups of pixels) of RS imagery. Considering semantic segmentation of 3D point clouds, it is to classify points of a 3D point cloud into categories. TDs are usually composed of RS images/point clouds, and the corresponding labelled value of each pixel/point recording its class.
- Change detection. These algorithms identify the difference between images acquired over the same geographical area but taken at different times. The TD comprise a set of pre-change and post-change RS images, with the corresponding ground truth map labelled changed and unchanged pixels. The image can be optical or SAR images.
- 3D model reconstruction. These algorithms infer the 3D geometry and structure of objects and scenes, mainly realized from the dense matching of multi-view images. The TD are usually composed of two-view or multi-view images, with the corresponding disparity map or depth maps as ground truth respectively.

6.2. Modularization

The TrainingDML-AI conceptual model provides models for the most important elements within TD. These elements have been identified to be either required or important in many different AI/ML tasks. However, implementations are not required to support the complete TrainingDML-AI model in order to be conformant to the Standard. Implementations may employ a subset of constructs according to their specific information needs. For this purpose, modularization is applied to the TrainingDML-AI.

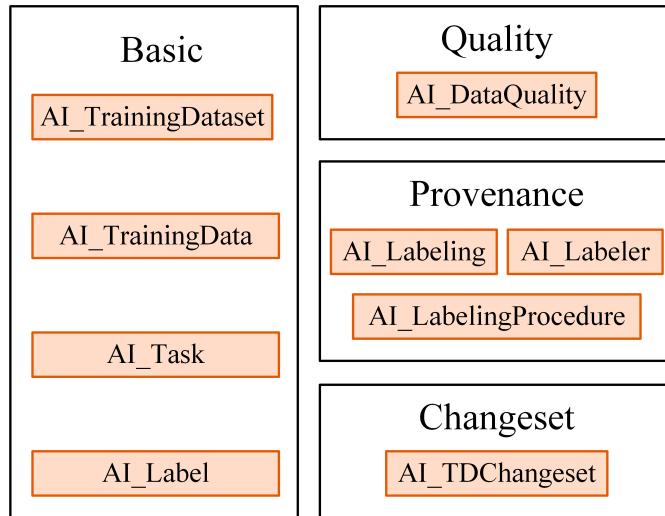


Figure 2. TrainingDML-AI module overview.

As shown in Figure 2, the TrainingDML-AI conceptual model is thematically decomposed into a Basic module, a Provenance module, a Quality module and a Changeset module. The Basic module comprises the basic concepts and elements, including AI_TrainingDataset, AI_TrainingData, AI_Label, and AI_Task, of the TrainingDML-AI, and thus, must be implemented by any conformant system. The Provenance module provides a comprehensive definition of provenance by AI_Labeling, AI_Labeler, and AI_Labeling Procedure. The Quality module offers quality description of TD with AI_DataQuality elements. And the Changeset module defines AI_TDChangeset between versions of datasets.

6.3. General modeling principles

6.3.1. Element modeling

The modeling of all elements in the TrainingDML-AI conceptual model has the following principles (Reference [7]):

- Granularity. Two levels of granularity are differentiated in the conceptual model: The Training Dataset is used to refer to the collection level, and the Training Data is used to refer to the individual level.
- Label semantics. The training dataset will not be limited to one classification scheme. External classification schemes should be allowed to be linked into the Training Dataset to accommodate different cases in practice.
- Light-weight design. The lightweight designed conceptual model has a minimum set of metadata elements, provenance, or quality measures at the collection level instead of at the

individual level. This is to facilitate the understanding of the dataset and improve the scalability for communicating large training datasets.

- Alignment. The modelling of elements in TDs can leverage existing efforts for wide adoption, such as for ISO 19109 Geographic information — Rules for application schema, ISO 19115-1 Geographic information — Metadata — Part 1: Fundamentals, ISO 19157-1 Geographic information — Data quality — Part 1: General requirements, and the OGC Geography Markup Language (GML) Standard. The conceptual model can be aligned with these existing standards and leverage capabilities fulfilled in part by other standards.
- Quality, bias, and ethics. Elements related to quality, or more specifically, bias that can be used to reduce the errors when using AI/ML. For example, any knowledge of the TD imbalance and mislabeling can be stored in TD quality. In addition, data ethics aims to safeguard the responsible use of TD, and it can be addressed by using the license property in the TD.
- Changeset. This will be an optional module in TD modelling. Changeset addresses how to capture changes in TD datasets. The change model considers the trend in TD collections to use the crowdsourcing platforms and borrow the change representation from the platforms such as OpenStreetMap.

6.3.2. Class Hierarchy and Inheritance of Properties and Relations

In the TrainingDML-AI conceptual model, the specific elements such as EO training datasets, EO training data, scene label, object label, and pixel label are defined as subclasses of more general higher-level classes. Hence, elements build a hierarchy along specialization / generalization relationships where more specialized elements inherit the properties and relationships of all their super classes along the entire generalization path to the topmost element.

6.3.3. Definition of the Semantics for all Classes, Properties, and Relations

The meanings of all elements defined in the TrainingDML-AI conceptual model are normatively specified in the data dictionary in Clause 8.

6.3.4. Data Integrity, Authenticity, and Non-repudiation

Sometimes training datasets can be downloaded, disseminated, and changed by anyone. The data integrity, authenticity, and non-repudiation are important to ensure unexpected bias propagation and distorted results. Currently the standard focuses on the information modelling, while data dissemination can be enriched with strategies from the general information domain by publishing hashes (e.g., MD5) and public-keys (e.g., RSA) after signing and encrypting.

6.4. Extending TrainingDML-AI

The TrainingDML-AI conceptual model is designed as a universal information model that defines elements and attributes which are useful for a broad range of AI/ML applications. In practical AI/ML applications, the elements within specific TDs will most likely contain attributes which are not explicitly modeled in TrainingDML-AI. Moreover, there might be TD elements which are not covered by the TrainingDML-AI thematic classes.

The model provides an abstract class-based method to support the exchange of such data. Elements

not represented by the predefined thematic classes of the model may be modeled and exchanged by extending abstract class.

Chapter 7. TrainingDML-AI UML Model

The TrainingDML-AI UML model is the normative definition of the TrainingDML-AI Conceptual Model. The tables and figures in this section were software generated from the UML model. As such, this section provides a normative representation of the TrainingDML-AI Conceptual Model.

7.1. ISO dependencies

TrainingDML-AI builds on the ISO 19100 family of standards. The applicable standards are identified in Figure 3. Data dictionaries are included for all the ISO-defined classes explicitly referenced in the TrainingDML-AI UML model. These data dictionaries are provided for the convenience of the user. The ISO standards are the normative source.

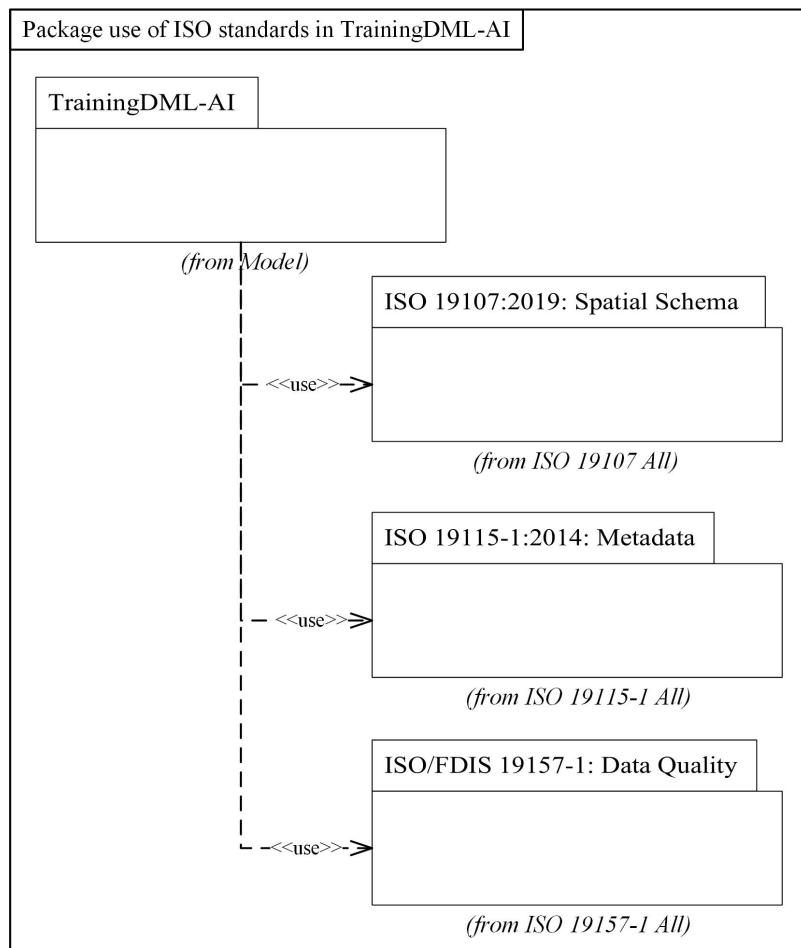


Figure 3. Use of ISO Standards in TrainingDML-AI.

The ISO classes explicitly used in the TrainingDML-AI UML model are introduced in Table 1. Further details about these classes can be found in the Data Dictionary within Clause 8.

Table 1. ISO Classes used in TrainingDML-AI

Class Name	Description
Feature	Abstraction of real world phenomena.
MD_Band	Range of wavelengths in the electromagnetic spectrum.

MD_Scope	The target resource and physical extent for which information is reported.
EX_Extent	Extent of the resource.
CI_Citation	Standardized resource reference.
DataQuality	Quality information for the data specified by a data quality scope.
QualityElement	Aspect of quantitative quality information.

7.2. Overview of the UML model

The UML model is presented with core concepts in Figure 4, followed by the concrete classes in Figure 5. The following describes the core concepts:

- AI_TrainingDataset: This concept represents a collection of training samples, i.e. a training dataset.
- AI_TrainingData: This concept is an individual training sample in a training dataset.
- AI_Task: This concept is used to identify the task that the training dataset is used for.
- AI_Label: This concept represents the label semantics for TD.
- AI_Labeling: This concept provides the provenance of how TD are created.
- AI_TDChangeset: This concept records types of TD changes between two versions of the training dataset.
- AI_DataQuality: This concept is associated with a training dataset to document its quality.

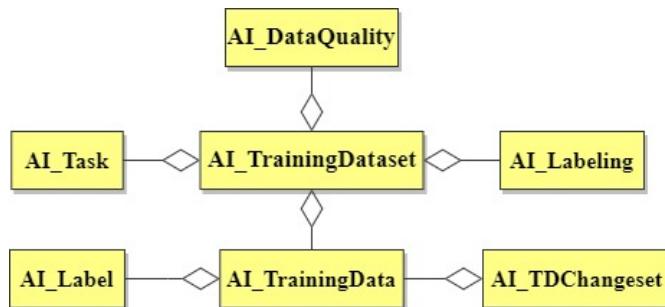


Figure 4. Core concepts.

The full overview of concrete classes and attributes are presented in Figure 5. Concepts related to the EO AI/ML applications are defined as classes extended from abstract classes. Each core concept with related classes will be described in the rest subsections.

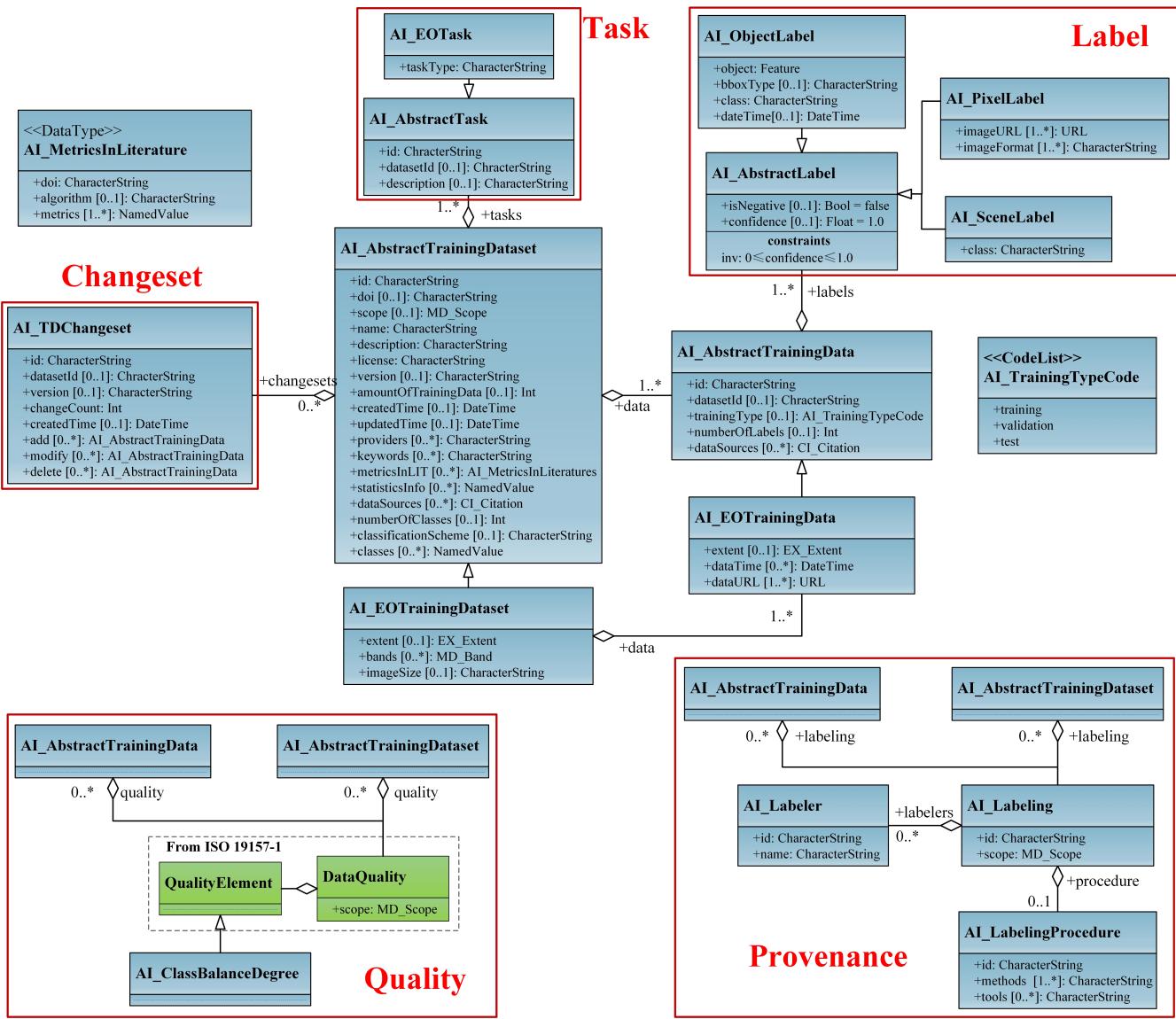


Figure 5. Overview of the UML model.

7.3. AI_TrainingDataset

Requirements class 1

<http://www.opengis.net/spec/TrainingDML-AI-1/1.0/req/aitrainingdataset>

Target type	Implementation Standard
Dependency	ISO 19115-1:2014
Dependency	/req/aitrainingdata
Dependency	/req/aitask
Dependency	/req/ailabeling
Dependency	/req/aidataquality
Dependency	/req/aitdchangeset

An AI training dataset is represented as **AI_AbstractTrainingDataset**. Each training sample in the dataset is represented as **AI_AbstractTrainingData**. A set of basic attributes are defined at the collection level for **AI_AbstractTrainingDataset**. These include the identification, authorship, time

for creation and update, AI tasks, number of training samples (amountOfTrainingData), label semantics (classes recording the available semantic classes), number of label classes, and statistics of training samples in each class (statisticsInfo). These basic attributes at the collection level can also be seen as a minimum set of metadata to define a simple payload. In this context, payload means the actual content that will be transmitted on the Web. They can be mapped to the MD_Metadata entity properties in ISO 19115, and thus delegate the metadata description via a metadata association to existing metadata standards.

The AI_EOTrainingDataset is defined to convey attributes specific to EO. For example, the data source (e.g., EO images), image size, band information, and spatial extent are defined in AI_EOTrainingDataset.

The UML diagram of the AI_TrainingDataset is illustrated in Figure 6.

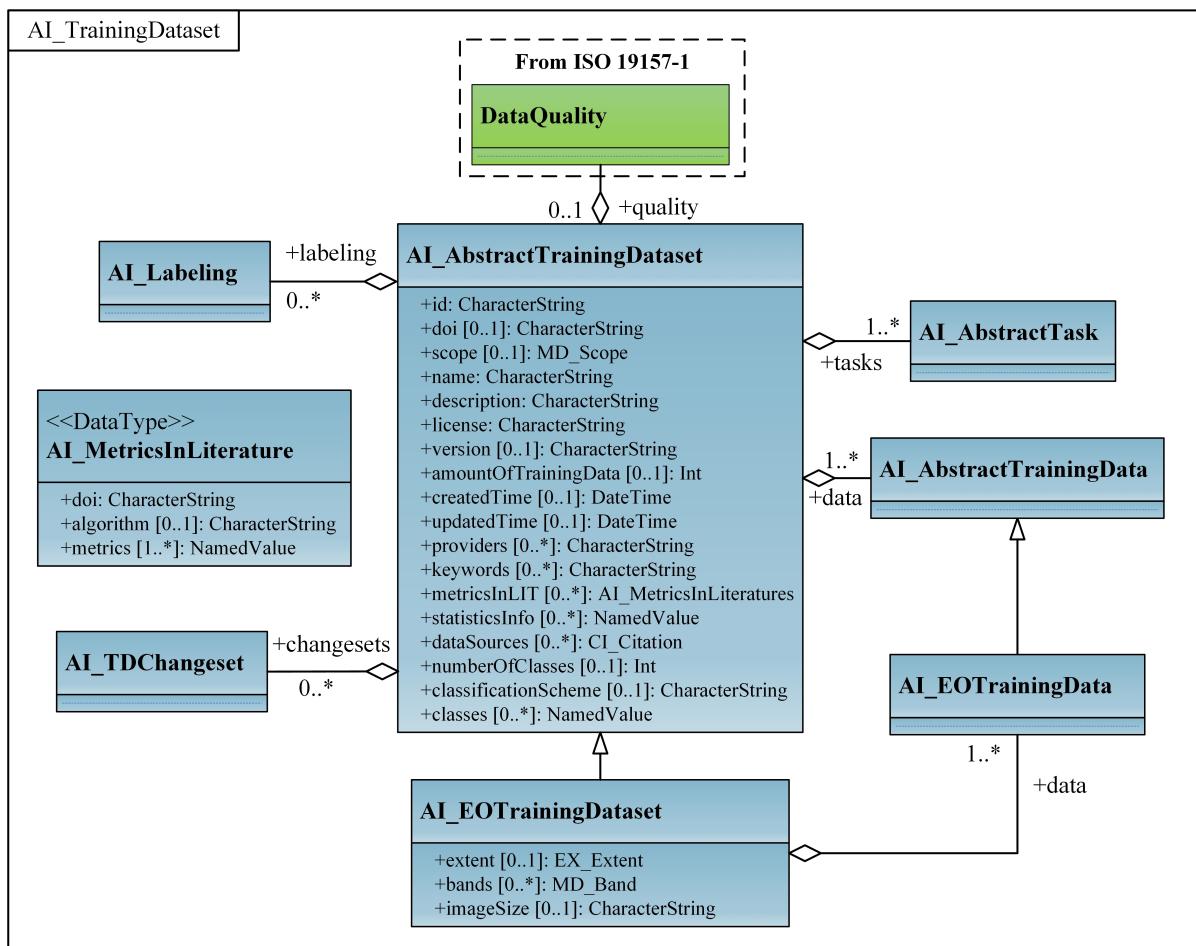


Figure 6. UML diagram of AI_TrainingDataset.

7.3.1. Requirements

The following requirement defines the rules governing the implementation of the AI_TrainingDataset conceptual model as an implementation standard.

Requirement 1

/req/airainingdataset/classes

For each UML class defined or referenced in AI_TrainingDataset conceptual model

Requirement 1	
A	Any implementation standard SHALL contain an element which represents the same concept as that defined for the UML class.
B	Any implementation standard SHALL represent associations with the same source, target, direction, roles, and multiplicities as those of the UML class.
C	Any implementation standard SHALL represent the attributes of the UML class including the name, definition, type, and multiplicity.
D	Any implementation standard SHALL represent the attributes of all super classes of the UML class including the name, definition, type, and multiplicity.
E	Any implementation standard SHALL represent the associations of all super classes of the UML class including the source, target, direction, roles, and multiplicity.
F	Any implementation standard SHALL specify how an implementation observes all constraints the Conceptual Model imposes on the UML class.

An implementing technology may not be able to support all the concepts defined in the TrainingDML-AI UML model. Alternately, some concepts may be inappropriate for the application domain for which the implementation standard was developed. In those cases, elements of the UML model may be mapped to null elements.

Permission 1	
/per/airtrainingdataset/classes	
For each UML class defined or referenced in AI_TrainingDataset	
A	An implementation standard MAY represent that class as a null class with no attributes, associations, or definition.
B	An implementation standard MAY represent an association of the UML class with a null association.
C	An implementation standard MAY represent an attribute of the UML class with a null attribute.
D	Whenever a null element is used to represent a concept from the AI_TrainingDataset, the implementation standard SHOULD document that mapping and provide an explanation for why that concept was not implemented.

7.3.2. Class definitions

Table 2. Classes defined in AI_TrainingDataset

Name	Description
AI_AbstractTrainingDataset	AI_AbstractTrainingDataset defines the basic concepts and components of different kinds of training dataset.
AI_EOTrainingDataset	AI_EOTrainingDataset describes attributes specific to EO training dataset. For example, the data source, image size, band information, and spatial extent.

Table 3. Data types defined in AI_TrainingDataset

Name	Description
AI_MetricsInLiterature <<DataType>>	AI_MetricsInLiterature records results of performance metrics achieved by AI/ML algorithms in the peer-reviewed literature.

7.4. AI_TrainingData

Requirements class 2

<http://www.opengis.net/spec/TrainingDML-AI-1/1.0/req/airtrainingdata>

Target type	Implementation Standard
Dependency	ISO 19115-1:2014
Dependency	/req/airtrainingdataset
Dependency	/req/ailabel

The AI_AbstractTrainingData includes the source data and its corresponding labels, which focus on a compact modelling of individual training samples. When there is a need to identify the different uses/purposes, an optional attribute, the training type (training, validation, or test types), can be added at the individual level. The attribute datasetId helps identify samples from different training datasets.

The AI_EOTrainingData is defined to convey attributes specific to EO. It provides additional attributes, including spatial extent and date/time of the source EO inputs.

The UML diagram of the AI_TrainingData is illustrated in Figure 7.

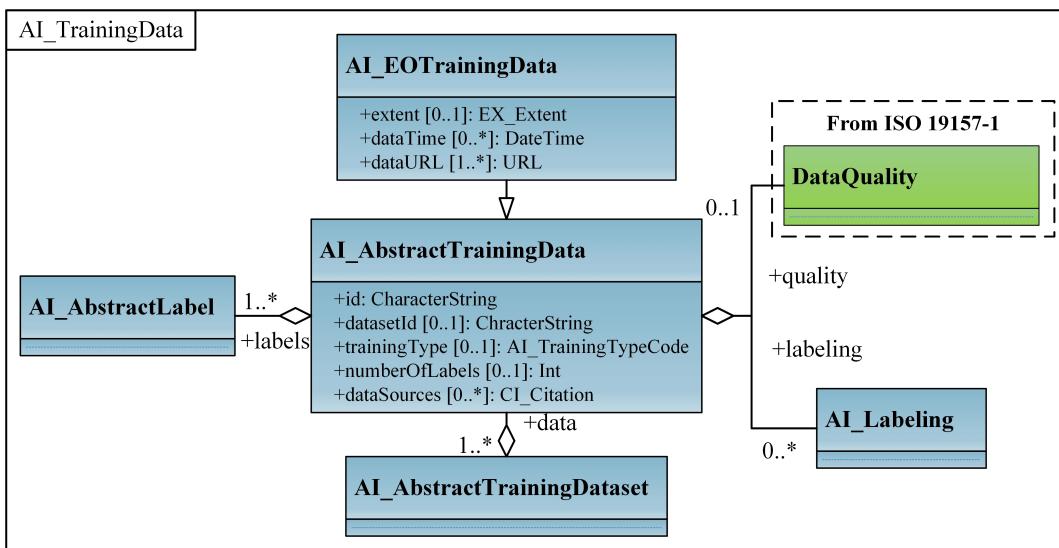


Figure 7. UML diagram of AI_TrainingData.

The Code Lists provided for the AI_TrainingData are illustrated in Figure 8.

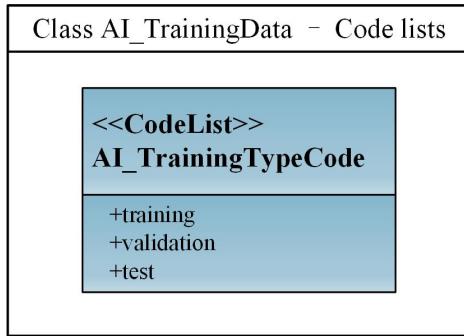


Figure 8. Code lists from the AI_TrainingData.

7.4.1. Requirements

The following requirement defines the rules governing the implementation of the AI_TrainingData conceptual model as an implementation standard.

Requirement 2

/req/aitrainingdata/classes

For each UML class defined or referenced in AI_TrainingData

A	Any implementation standard SHALL contain an element which represents the same concept as that defined for the UML class.
B	Any implementation standard SHALL represent associations with the same source, target, direction, roles, and multiplicities as those of the UML class.
C	Any implementation standard SHALL represent the attributes of the UML class, including the name, definition, type, and multiplicity.
D	Any implementation standard SHALL represent the attributes of all super classes of the UML class, including the name, definition, type, and multiplicity.
E	Any implementation standard SHALL represent the associations of all super classes of the UML class, including the source, target, direction, roles, and multiplicity.
F	Any implementation standard SHALL specify how an implementation observes all constraints the Conceptual Model imposes on the UML class.

Permission 2

/per/aitrainingdata/classes

For each UML class defined or referenced in AI_TrainingData

A	An implementation standard MAY represent that class as a null class with no attributes, associations, or definition.
B	An implementation standard MAY represent an association of the UML class with a null association.
C	An implementation standard MAY represent an attribute of the UML class with a null attribute.

Permission 2	
D	Whenever a null element is used to represent a concept from the AI_TrainingData, the implementation standard SHOULD document that mapping and provide an explanation for why that concept was not implemented.

7.4.2. Class definitions

Table 4. Classes defined in AI_TrainingData

Name	Description
AI_AbstractTrainingData	AI_AbstractTrainingData defines the basic concepts and components of different kinds of training data.
AI_EOTrainingData	AI_EOTrainingData describes attributes specific to EO training data. For example, spatial extent and date time of the source EO inputs.

Table 5. Code list classes defined in AI_TrainingData

Name	Description
AI_TrainingTypeCode <<CodeList>>	AI_TrainingTypeCode is a code list used to identify whether the individual data element is used for different purposes.

7.5. AI_Task

Requirements class 3	
http://www.opengis.net/spec/TrainingDML-AI-1/1.0/req/aitask	
Target type	Implementation Standard
Dependency	/req/airainingdataset

Various AI/ML tasks in specific domains can have different organizing forms of TD, including source data and label representations. AI_EOTask is proposed by extending AI_AbstractTask to represent specific AI/ML tasks in the EO domain. The task type can refer to a particular type defined by an external category, such as scene classification or change detection. It specifies, for example, whether a remote sensing training dataset is used for scene classification, objection detection, land use or land cover, or change detection scenarios.

The UML diagram of the AI_Task is illustrated in Figure 9.

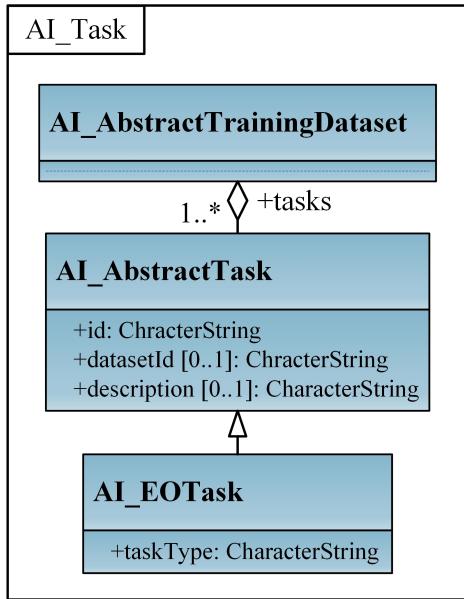


Figure 9. UML diagram of *AI_Task*.

7.5.1. Requirements

The following requirement defines the rules governing the implementation of the *AI_Task* conceptual model as an implementation standard.

Requirement 3

/req/aitask/classes

For each UML class defined or referenced in *AI_Task*

A	Any implementation standard SHALL contain an element which represents the same concept as that defined for the UML class.
B	Any implementation standard SHALL represent associations with the same source, target, direction, roles, and multiplicities as those of the UML class.
C	Any implementation standard SHALL represent the attributes of the UML class including the name, definition, type, and multiplicity.
D	Any implementation standard SHALL represent the attributes of all super classes of the UML class including the name, definition, type, and multiplicity.
E	Any implementation standard SHALL represent the associations of all super classes of the UML class including the source, target, direction, roles, and multiplicity.
F	Any implementation standard SHALL specify how an implementation observes all constraints the Conceptual Model imposes on the UML class.

Permission 3

/per/aitask/classes

For each UML class defined or referenced in *AI_Task*

A	An implementation standard MAY represent that class as a null class with no attributes, associations, or definition.
---	--

Permission 3	
B	An implementation standard MAY represent an association of the UML class with a null association.
C	An implementation standard MAY represent an attribute of the UML class with a null attribute.
D	Whenever a null element is used to represent a concept from the AI_Task, the implementation standard SHOULD document that mapping and provide an explanation for why that concept was not implemented.

7.5.2. Class definitions

Table 6. Classes defined in AI_Task

Name	Description
AI_AbstractTask	AI_AbstractTask defines the AI/ML tasks in specific domains that have different organizing forms of TD.
AI_EOTask	AI_EOTask extends AI_AbstractTask to represent specific AI/ML tasks in the EO domain.

7.6. AI_Label

Requirements class 4

<http://www.opengis.net/spec/TrainingDML-AI-1/1.0/req/ailabel>

Target type	Implementation Standard
Dependency	/req/airainingdata

Labels for each individual training sample can be represented using a feature, coverage, or a semantic class from ontologies or shared vocabularies (external classification schemes and classes). A training sample typically relates pixels, objects, and scenes to semantic labels, where labels are encoded as coverage, geometric polygon, and text respectively. Therefore, the AI_AbstractLabel is extended to specify AI_SceneLabel, AI_ObjectLabel, and AI_PixelLabel respectively. There is no restriction for TD producers to describe label information with label classes. For example, labels of training dataset for object detection can be represented by either AI_ObjectLabel or AI_PixelLabel.

The UML diagram of the AI_Label is illustrated in Figure 10.

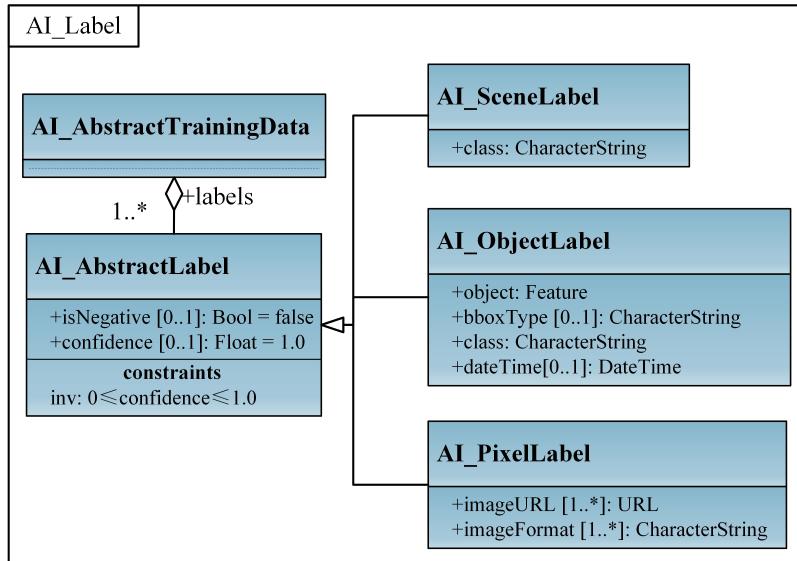


Figure 10. UML diagram of AI_Label.

7.6.1. Requirements

The following requirement defines the rules governing the implementation of the AI_Label conceptual model as an implementation standard.

Requirement 4

/req/ailabel/classes

For each UML class defined or referenced in AI_Label

A	Any implementation standard SHALL contain an element which represents the same concept as that defined for the UML class.
B	Any implementation standard SHALL represent associations with the same source, target, direction, roles, and multiplicities as those of the UML class.
C	Any implementation standard SHALL represent the attributes of the UML class including the name, definition, type, and multiplicity.
D	Any implementation standard SHALL represent the attributes of all super classes of the UML class including the name, definition, type, and multiplicity.
E	Any implementation standard SHALL represent the associations of all super classes of the UML class including the source, target, direction, roles, and multiplicity.
F	Any implementation standard SHALL specify how an implementation observes all constraints the Conceptual Model imposes on the UML class.

Permission 4

/per/ailabel/classes

For each UML class defined or referenced in AI_Label

A	An implementation standard MAY represent that class as a null class with no attributes, associations, or definition.
---	--

Permission 4	
B	An implementation standard MAY represent an association of the UML class with a null association.
C	An implementation standard MAY represent an attribute of the UML class with a null attribute.
D	Whenever a null element is used to represent a concept from the AI_Label, the implementation standard SHOULD document that mapping and provide an explanation for why that concept was not implemented.

7.6.2. Class definitions

Table 7. Classes defined in AI_Label

Name	Description
AI_AbstractLabel	AI_AbstractLabel defines a set of informative tags or metadata obtained by labelling training data. It can be represented using a feature, coverage, or a semantic class from ontologies or shared vocabularies.
AI_SceneLabel	AI_SceneLabel represents the scene level label using a semantic from ontologies or shared vocabularies.
AI_ObjectLabel	AI_ObjectLabel represents the object level label using a feature.
AI_PixelLabel	AI_PixelLabel represents the pixel level label using a coverage.

7.7. AI_Labeling

Requirements class 5

<http://www.opengis.net/spec/TrainingDML-AI-1/1.0/req/ailabeling>

Target type	Implementation Standard
Dependency	/req/airainingdataset

AI_Labeling can be associated with AI_AbstractTrainingDataset or AI_AbstractTrainingData to record basic provenance information on how the training dataset or training samples are created. For example, AI_LabelingProcedure can be used to record the sampling strategy, and dataSources from AI_AbstractTrainingDataset/ AI_AbstractTrainingData can refer to any additional geospatial data that was used as part of the sampling strategy.

AI_Labeling includes the labeler and the labeling procedure, which can be mapped to the agent and activity respectively in W3C PROV model. The labeler identifies the agent that creates the training dataset or individual samples, and the labeling procedure represents the process for data generation. AI_Labeling can also describe how the input data of the TD has been manipulated, such as resampled, color corrected, atmospherically corrected, and terrain corrected. For example, imagery with 10cm and 15cm resolution was resampled to 20cm prior for labeling.

Alternatively, despite the simple payload in this TrainingDML-AI Standard, providing a

comprehensive definition of provenance using the ISO 19115 lineage model through a metadata association in the data is also applicable.

The UML diagram of the AI_Labeling is illustrated in Figure 11.

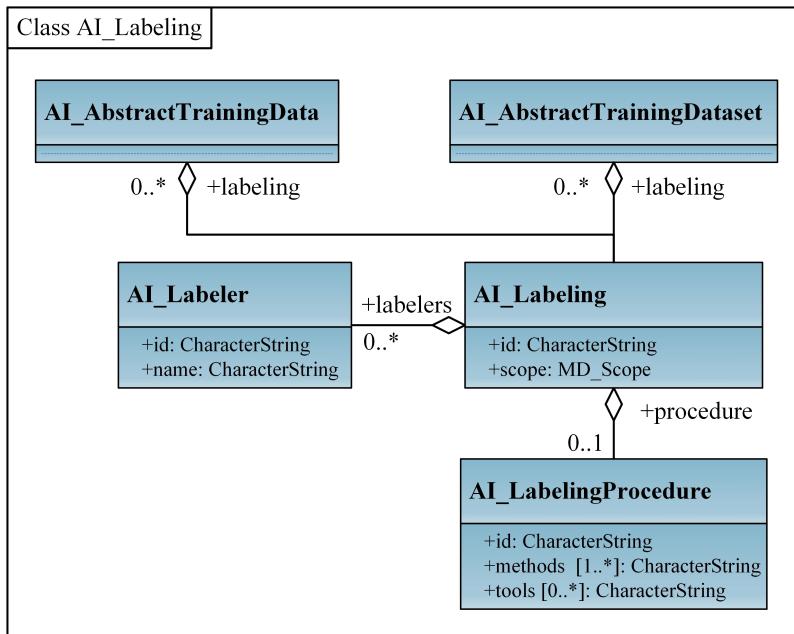


Figure 11. UML diagram of AI_Labeling.

7.7.1. Requirements

The following requirement defines the rules governing the implementation of the AI_Labeling conceptual model as an implementation standard.

Requirement 5	
/req/ailabeling/classes	
For each UML class defined or referenced in AI_Labeling	
A	Any implementation standard SHALL contain an element which represents the same concept as that defined for the UML class.
B	Any implementation standard SHALL represent associations with the same source, target, direction, roles, and multiplicities as those of the UML class.
C	Any implementation standard SHALL represent the attributes of the UML class including the name, definition, type, and multiplicity.
D	Any implementation standard SHALL represent the attributes of all super classes of the UML class, including the name, definition, type, and multiplicity.
E	Any implementation standard SHALL represent the associations of all super classes of the UML class, including the source, target, direction, roles, and multiplicity.
F	Any implementation standard SHALL specify how an implementation observes all constraints the Conceptual Model imposes on the UML class.

Permission 5	
/per/ailabeling/classes	
For each UML class defined or referenced in AI_Labeling	
A	An implementation standard MAY represent that class as a null class with no attributes, associations, or definition.
B	An implementation standard MAY represent an association of the UML class with a null association.
C	An implementation standard MAY represent an attribute of the UML class with a null attribute.
D	Whenever a null element is used to represent a concept from the AI_Labeling, the implementation standard SHOULD document that mapping and provide an explanation for why that concept was not implemented.

7.7.2. Class definitions

Table 8. Classes defined in AI_Labeling

Name	Description
AI_Labeling	AI_Labeling defines basic provenance information on how to create the training dataset and includes the labeler and labeling procedure.
AI_Labeler	AI_Labeler identifies the agent that creates the labels and can be mapped to the agent in W3C PROV.
AI_LabelingProcedure	AI_LabelingProcedure represents the process for labeling and can be mapped to the activity in W3C PROV.

7.8. AI_TDChangeset

Requirements class 6

<http://www.opengis.net/spec/TrainingDML-AI-1/1.0/req/aitdchangeset>

Target type	Implementation Standard
Dependency	/req/airainingdata

AI_TDChangeset records changed training samples between two versions in the collection level, including added training samples, modified training samples and deleted training samples. AI_TDChangeset includes three kinds of changes (add, modify, and delete) for individuals.

The UML diagram of the AI_TDChangeset is illustrated in Figure 12.

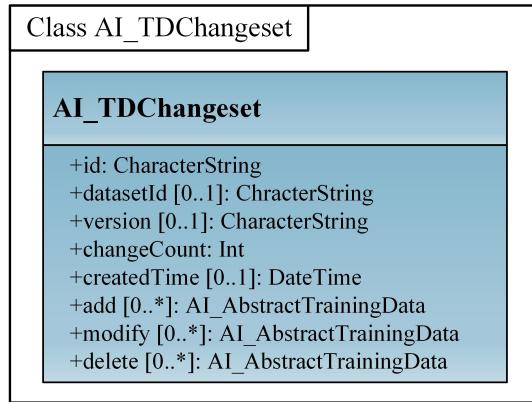


Figure 12. UML diagram of *AI_TDChangeset*

7.8.1. Requirements

The following requirement defines the rules governing the implementation of the *AI_TDChangeset* conceptual model as an implementation standard.

Requirement 6	
/req/aitdchangeset/classes	
For each UML class defined or referenced in <i>AI_TDChangeset</i>	
A	Any implementation standard SHALL contain an element which represents the same concept as that defined for the UML class.
B	Any implementation standard SHALL represent associations with the same source, target, direction, roles, and multiplicities as those of the UML class.
C	Any implementation standard SHALL represent the attributes of the UML class, including the name, definition, type, and multiplicity.
D	Any implementation standard SHALL represent the attributes of all super classes of the UML class, including the name, definition, type, and multiplicity.
E	Any implementation standard SHALL represent the associations of all super classes of the UML class, including the source, target, direction, roles, and multiplicity.
F	Any implementation standard SHALL specify how an implementation observes all constraints the Conceptual Model imposes on the UML class.

Permission 6	
/per/aitdchangeset/classes	
For each UML class defined or referenced in <i>AI_TDChangeset</i>	
A	An implementation standard MAY represent that class as a null class with no attributes, associations, or definition.
B	An implementation standard MAY represent an association of the UML class with a null association.
C	An implementation standard MAY represent an attribute of the UML class with a null attribute.

Permission 6	
D	Whenever a null element is used to represent a concept from the AI_TDChangeset, the implementation standard SHOULD document that mapping and provide an explanation for why that concept was not implemented.

7.8.2. Class definitions

Table 9. Classes defined in AI_TDChangeset

Name	Description
AI_TDChangeset	AI_TDChangeset represents a set of individual level changes (add, modify, and delete) of the training dataset.

7.9. AI_DataQuality

Requirements class 7

<http://www.opengis.net/spec/TrainingDML-AI-1/1.0/req/aidataquality>

Target type	Implementation Standard
Dependency	ISO 19157-1
Dependency	/req/aitrainingdataset

TD quality description can use DataQuality from ISO 19157-1 to align with the existing efforts on geographic data quality. Data quality can be evaluated in terms of either collection or individual levels of the TD. Although most data quality elements are designed for datasets, it is sometimes needed to know the data quality at the individual level. For example, how is the positional uncertainty of an individual training sample measured by a GPS device in the field? Also in terms of the object label, an individual training sample includes a source image and many object labels. In this case, the quality can be measured based on these labels, such as Completeness. Thus the scope of the data quality can be used to specify the level and extent that identify the data on which data quality is to be established and evaluated. For TrainingDML-AI, AI_TrainingDataset/AI_TrainingData can be mapped to dataset/feature in MD_ScopeCode.

The quality description of the TD can also leverage the quality elements defined in QualityElement. For example, Completeness describes label commission and omission, Thematic accuracy describes class accuracy of labels. TD related quality elements like AI_ClassBalanceDegree can be defined by extending the QualityElement. AI_ClassBalanceDegree can help address the bias issue and evaluate whether the number of classes in the training dataset is imbalanced.

The UML diagram of the AI_DataQuality module is illustrated in Figure 13.

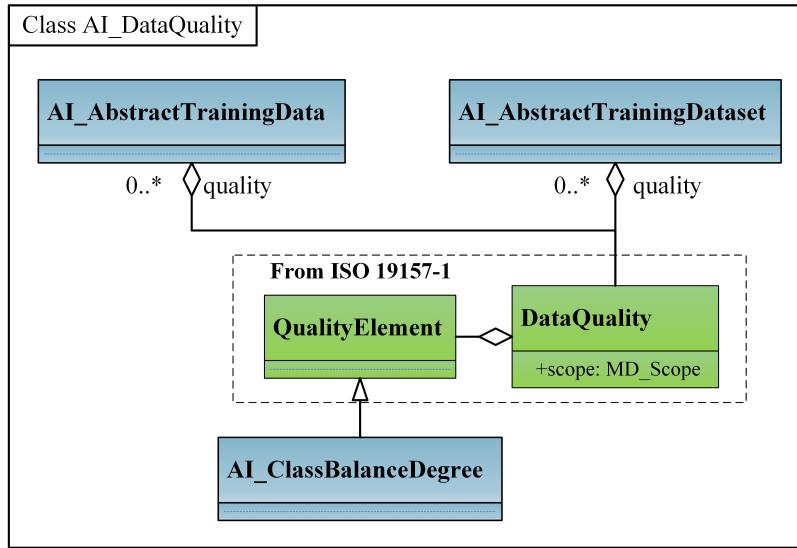


Figure 13. UML diagram of `AI_DataQuality`.

7.9.1. Requirements

The following requirement defines the rules governing the implementation of the `DataQuality` conceptual model as an implementation standard.

Requirement 7	
/req/aidataquality/classes	
For each UML class defined or referenced in <code>DataQuality</code>	
A	Any implementation standard SHALL contain an element which represents the same concept as that defined for the UML class.
B	Any implementation standard SHALL represent associations with the same source, target, direction, roles, and multiplicities as those of the UML class.
C	Any implementation standard SHALL represent the attributes of the UML class, including the name, definition, type, and multiplicity.
D	Any implementation standard SHALL represent the attributes of all super classes of the UML class, including the name, definition, type, and multiplicity.
E	Any implementation standard SHALL represent the associations of all super classes of the UML class, including the source, target, direction, roles, and multiplicity.
F	Any implementation standard SHALL specify how an implementation observes all constraints the Conceptual Model imposes on the UML class.

Permission 7	
/per/aidataquality/classes	
For each UML class defined or referenced in <code>DataQuality</code>	
A	An implementation standard MAY represent that class as a null class with no attributes, associations, or definition.
B	An implementation standard MAY represent an association of the UML class with a null association.

Permission 7	
C	An implementation standard MAY represent an attribute of the UML class with a null attribute.
D	Whenever a null element is used to represent a concept from the DataQuality, the implementation standard SHOULD document that mapping and provide an explanation for why that concept was not implemented.

7.9.2. Class definitions

Table 10. Classes defined in AI_DataQuality

Name	Description
AI_ClassBalanceDegree	AI_ClassBalanceDegree defines the quality elements for measuring degree of class balance of training datasets.

Chapter 8. TrainingDML-AI Data Dictionary

The TrainingDML-AI UML model is the normative definition of the TrainingDML-AI Conceptual Model. The Data Dictionary tables in this section were software generated from the UML model. As such, this section provides a normative representation of the TrainingDML-AI Conceptual Model.

8.1. ISO Classes

The following classes are defined in the ISO standards and used by the TrainingDML-AI Conceptual Model.

8.1.1. Feature (from ISO 19107:2019)

Table 11. Metadata of Feature (Class)

Definition	Abstraction of real world phenomena
Subclass of	None
Stereotype	<<Type>>

8.1.2. MD_Band (from ISO 19115-1:2014)

Table 12. Metadata of MD_Band (Class)

Definition	MD_Band is the range of wavelengths in the electromagnetic spectrum.
Subclass of	MD_SampleDimension
Stereotype	<<Class>>

8.1.3. MD_Scope (from ISO 19115-1:2014)

Table 13. Metadata of MD_Scope (DataType)

Definition	The target resource and physical extent for which information is reported.
Subclass of	None
Stereotype	<<DataType>>

8.1.4. EX_Extent (from ISO 19115-1:2014)

Table 14. Metadata of EX_Extent (DataType)

Definition	EX_Extent is the extent of the resource electromagnetic spectrum.
Subclass of	None
Stereotype	<<DataType>>

8.1.5. CI_Citation (from ISO 19115-1:2014)

Table 15. Metadata of CI_Citation (Class)

Definition	CI_Citation is the standardized resource reference.
Subclass of	None
Stereotype	<<Class>>

8.1.6. DataQuality (from ISO 19157-1)

Table 16. Metadata of DataQuality (Class)

Definition	DataQuality is the quality information for the data specified by a data quality scope.
Subclass of	None
Stereotype	<<Class>>

8.1.7. QualityElement (from ISO 19157-1)

Table 17. Metadata of QualityElement (Class)

Definition	QualityElement is the aspect of quantitative quality information.
Subclass of	None
Stereotype	<<Class>>

8.2. AI_TrainingDataset

Table 18. Metadata of AI_TrainingDataset (ApplicationSchema)

Description	The AI_TrainingDataset module supports the modelling of different kinds of training dataset.
Parent package	TrainingDML-AI
Stereotype	<<ApplicationSchema>>

8.2.1. Classes

AI_AbstractTrainingDataset

Table 19. Metadata of AI_AbstractTrainingDataset (Class)

Definition	AI_AbstractTrainingDataset defines the basic concepts and components of different kinds of training dataset.
Subclass of	None
Stereotype	<<Class>>

Table 20. Attributes of AI_AbstractTrainingDataset (Class)

Attribute	Value type and multiplicity	Definition
id	CharacterString [1..1]	Identification of the AI training dataset.
doi	CharacterString [0..1]	Digital object identifier of the AI training dataset.
scope	MD_Scope [0..1]	Description of the scope of the training dataset.
name	CharacterString [1..1]	Name of the AI training dataset.
description	CharacterString [1..1]	Description of the AI training dataset.
version	CharacterString [0..1]	Version number of the AI training dataset.
amountOfTrainingData	Int [1..1]	Total number of training samples in the AI training dataset.
createdTime	DateTime [0..1]	Time when the AI training dataset was created.
updatedTime	DateTime [0..1]	Time when the AI training dataset was updated.
license	CharacterString [0..1]	License description of the AI training dataset.
providers	CharacterString [0..*]	People or organizations who provide the AI training dataset.
keywords	CharacterString [0..*]	Keywords of the AI training dataset.
metricsInLIT	AI_MetricsInLiterature [0..*]	Results of performance metrics achieved by AI/ML algorithms in the peer-reviewed literature.
statisticsInfo	NamedValue [0..*]	Statistics results of training samples in each class.
dataSources	CI_Citation [0..*]	Citation of data sources.
numberOfClasses	Int [1..1]	Total number of classes in the AI training dataset.
classificationSchema	CharacterString [0..1]	Classification schema for classes used in the AI training dataset.
classes	NamedValue [1..1]	Classes used in the AI training dataset.

AI_EOTrainingDataset

Table 21. Metadata of AI_EOTrainingDataset (Class)

Definition	AI_EOTrainingDataset describes attributes specific to EO training dataset. For example, the data source, image size, band information, and spatial extent.
Subclass of	AI_AbstractTrainingDataset

Stereotype	<<Class>>
------------	-----------

Table 22. Attributes of AI_EOTrainingDataset (Class)

Attribute	Value type and multiplicity	Definition
extent	EX_Extent [0..1]	Spatial extent of the EO training dataset.
bands	MD_Bands [0..*]	Description of the image bands used in the EO training dataset.
imageSize	CharacterString [0..1]	Size of the images used in the EO training dataset.

AI_MetricsInLiterature

Table 23. Metadata of AI_MetricsInLiterature (DataType)

Definition	AI_MetricsInLiterature records results of performance metrics achieved by AI/ML algorithms in the peer-reviewed literature.
Subclass of	None
Stereotype	<<DataType>>

Table 24. Attributes of AI_MetricsInLiterature (DataType)

Attribute	Value type and multiplicity	Definition
doi	CharacterString [1..1]	Digital object identifier of the peer-reviewed literature.
algorithm	CharacterString [0..1]	AI/ML algorithms used in the peer-reviewed literature.
metrics	NamedValue [1..*]	Metrics and results of AI/ML algorithms in the peer-reviewed literature.

8.3. AI_TrainingData

Table 25. Metadata of AI_TrainingData (ApplicationSchema)

Description	AI_TrainingData module supports the modelling of an individual AI training sample.
Parent package	TrainingDML-AI
Stereotype	<<ApplicationSchema>>

8.3.1. Classes

AI_AbstractTrainingData

Table 26. Metadata of AI_AbstractTrainingData (Class)

Definition	AI_AbstractTrainingData defines the basic concepts and components of different kinds of training data.
Subclass of	None
Stereotype	<<Class>>

Table 27. Attributes of AI_AbstractTrainingData (Class)

Attribute	Value type and multiplicity	Definition
id	CharacterString [1..1]	Identification of the individual AI training sample.
datasetId	CharacterString [0..1]	Identification of the training dataset that the training sample belongs to.
trainingType	AI_TrainingTypeCode [0..1]	Training type of the individual AI training sample.
numberOfLabels	Int [0..1]	Total number of labels in the individual AI training sample.
dataSources	CI_Citation [0..*]	Citation of inputs to prepare a training sample.

AI_EOTrainingData

Table 28. Metadata of AI_EOTrainingData (Class)

Definition	AI_EOTrainingData describes attributes specific to EO training data. For example, spatial extent and date time of the source EO inputs.
Subclass of	AI_AbstractTrainingData
Stereotype	<<Class>>

Table 29. Attributes of AI_EOTrainingData (Class)

Attribute	Value type and multiplicity	Definition
extent	EX_Extent [0..1]	Spatial extent of the individual EO training sample.
dateTime	DateTime [0..*]	Data times when the EO data were obtained.
dataURL	URL [1..*]	URLs of the EO data.

AI_TrainingTypeCode

Table 30. Metadata of AI_TrainingTypeCode (CodeList)

Definition	AI_TrainingTypeCode is a code list used to identify whether the individual is used for different purposes.
Subclass of	None
Stereotype	<<CodeList>>

8.4. AI_Task

Table 31. Metadata of AI_Task (ApplicationSchema)

Description	AI_Task module supports the modelling of different kinds of AI task.
Parent package	TrainingDML-AI
Stereotype	<>ApplicationSchema>>

8.4.1. Classes

AI_AbstractTask

Table 32. Metadata of AI_AbstractTask (Class)

Definition	AI_AbstractTask defines the AI/ML tasks in specific domains that have different organizing forms of TD.
Subclass of	None
Stereotype	<>Class>>

Table 33. Attributes of AI_AbstractTask (Class)

Attribute	Value type and multiplicity	Definition
id	CharacterString [1..1]	Identification of the task.
datasetId	CharacterString [0..1]	Identification of the training dataset the training sample belongs to.
description	CharacterString [0..1]	Description of the AI task.

AI_EOTask

Table 34. Metadata of AI_EOTask (Class)

Definition	AI_EOTask extends AI_AbstractTask to represent specific AI/ML tasks in the EO domain.
Subclass of	AI_AbstractTask
Stereotype	<>Class>>

Table 35. Attributes of AI_EOTask (Class)

Attribute	Value type and multiplicity	Definition
taskType	CharacterString [1..1]	Type description of the EO task.

8.5. AI_Label

Table 36. Metadata of AI_Label (ApplicationSchema)

Description	AI_Label module supports the modelling of different kinds of AI Label.
-------------	--

Parent package	TrainingDML-AI
Stereotype	<<ApplicationSchema>>

8.5.1. Classes

AI_AbstractLabel

Table 37. Metadata of AI_AbstractLabel (Class)

Definition	AI_AbstractLabel defines a set of informative tags or metadata obtained by labelling training data. AI_AbstractLabel can be represented using a feature, coverage, or a semantic class from ontologies or shared vocabularies.
Subclass of	None
Stereotype	<<Class>>
Constraint	valueBetween0and1: inv: $0 \leq \text{confidence} \leq 1.0$

Table 38. Attributes of AI_AbstractLabel (Class)

Attribute	Value type and multiplicity	Definition
isNegative	Bool [0..1]	Whether the training sample related to the label is a positive or negative sample.
confidence	Float [0..1]	Confidence score of the labeler.

AI_SceneLabel

Table 39. Metadata of AI_SceneLabel (Class)

Definition	AI_SceneLabel represents the scene level label using a semantic type from ontologies or shared vocabularies.
Subclass of	AI_AbstractLabel
Stereotype	<<Class>>

Table 40. Attributes of AI_SceneLabel (Class)

Attribute	Value type and multiplicity	Definition
class	CharacterString [1..1]	Class name that records the semantic type of the scene of the training sample.

AI_ObjectLabel

Table 41. Metadata of AI_ObjectLabel (Class)

Definition	AI_ObjectLabel represents the object level label using a feature.
Subclass of	AI_AbstractLabel
Stereotype	<<Class>>

Table 42. Attributes of AI_Label (Class)

Attribute	Value type and multiplicity	Definition
object	Feature [1..1]	Feature that represents the position and attributes of the object.
bboxType	CharacterString [0..1]	Type of the bbox.
class	CharacterString [1..1]	Class that records the semantic of the object type.
dateTime	DateTime [0..1]	Created time of the object label.

AI_PixelLabel

Table 43. Metadata of AI_PixelLabel (Class)

Definition	AI_PixelLabel represents the pixel level label using a coverage.
Subclass of	AI_AbstractLabel
Stereotype	<<Class>>

Table 44. Attributes of AI_PixelLabel (Class)

Attribute	Value type and multiplicity	Definition
imageURL	URL [1..*]	URL of the images representing the label information.
imageFormat	CharacterString [1..1]	Image data format.

8.6. AI_Labeling

Table 45. Metadata of AI_Labeling (ApplicationSchema)

Description	AI_Labeling module supports the modelling of provenance information of the training dataset.
Parent package	TrainingDML-AI
Stereotype	<<ApplicationSchema>>

8.6.1. Classes

AI_Labeling

Table 55. Metadata of AI_Labeling (Class)

Definition	AI_Labeling defines basic provenance information on how to create the training dataset. AI_Labeling includes the labeler and labeling procedure.
Subclass of	None
Stereotype	<<Class>>

Table 46. Attributes of Labeling (Class)

Attribute	Value type and multiplicity	Definition
id	CharacterString [1..1]	Identifier of the labeling.

AI_Labeler

Table 47. Metadata of AI_Labeler (Class)

Definition	AI_Labeler identifies the agent that creates the labels and can be mapped to the agent in W3C PROV.
Subclass of	None
Stereotype	<>Class>>

Table 48. Attributes of AI_Labeler (Class)

Attribute	Value type and multiplicity	Definition
id	CharacterString [1..1]	Identifier of the labeler.
name	CharacterString [1..1]	Name of the labeler.

AI_LabelingProcedure

Table 49. Metadata of AI_LabelingProcedure (Class)

Definition	AI_LabelingProcedure represents the process for labeling and can be mapped to the activity in W3C PROV.
Subclass of	None
Stereotype	<>Class>>

Table 50. Attributes of AI_LabelingProcedure (Class)

Attribute	Value type and multiplicity	Definition
id	CharacterString [1..1]	Identifier of the labeling procedure
methods	CharacterString [1..*]	Methods used in the labeling procedure.
tools	CharacterString [0..1]	Tools or software used in the labeling procedure.

8.7. AI_TDChangeset

Table 51. Metadata of AI_TDChangeset (ApplicationSchema)

Description	AI_TDChangeset module supports the modelling of change information of the training dataset.
Parent package	TrainingDML-AI
Stereotype	<>ApplicationSchema>>

8.7.1. Classes

AI_TDChangeset

Table 52. Metadata of AI_TDChangeset (Class)

Definition	AI_TDChangeset
Subclass of	None
Stereotype	<<Class>>

Table 53. Attributes of AI_TDChangeset (Class)

Attribute	Value type and multiplicity	Definition
id	CharacterString [1..1]	Identifier of the changeset.
datasetId	CharacterString [0..1]	Identifier of the training dataset the changeset belongs to.
version	CharacterString [0..1]	Version of the training dataset that the changeset belongs to.
changeCount	Int [1..1]	Total number of changed training samples.
createdTime	DateTime [0..1]	Created time of the changeset.
add	AI_AbstractTrainingData [0..*]	Added training samples.
modify	AI_AbstractTrainingData [0..*]	Modified training samples.
delete	AI_AbstractTrainingData [0..*]	Deleted training samples.

8.8. AI_DataQuality

Table 54. Metadata of AI_DataQuality (ApplicationSchema)

Description	AI_DataQuality module supports the modelling of quality description of different kinds of training datasets.
Parent package	TrainingDML-AI
Stereotype	<<ApplicationSchema>>

8.8.1. Classes

AI_ClassBalanceDegree

Table 55. Metadata of AI_ClassBalanceDegree (Class)

Definition	AI_ClassBalanceDegree defines the quality elements for measuring degree of class balance of training datasets.
Subclass of	QualityElement
Stereotype	<<Class>>

Table 56. Attributes of AI_ClassBalanceDegree (Class)

Attribute	Value type and multiplicity	Definition
measure	MeasureReference [1..1]	Reference to measure used
evaluationMethod	EvaluationMethod [1..1]	Evaluation information
result	QualityResult [1..*]	Value obtained from applying a data quality measure

Appendix A: Abstract Test Suite (Normative)

A.1. Introduction

The test method for TrainingDML-AI conceptual model specified in this ATS is manual inspection. Automated methods may be used where they exist.

A.2. Conformance class AI_TrainingDataset

Conformance class 1

<http://www.opengis.net/spec/TrainingDML-AI-1/1.0/conf/aitrainingdataset>

Requirements class	/req/aitrainingdataset
--------------------	------------------------

Abstract test A.1

/conf/aitrainingdataset/classes

Requirement	/req/aitrainingdataset/classes
-------------	--------------------------------

Test purpose	To validate that the implementation standard correctly implements the UML Classes defined in the Conceptual Model.
--------------	--

Test method type	Manual Inspection
------------------	-------------------

Abstract test A.1

Test method	<p>a. For each UML class defined or referenced in the AI_TrainingDataset Package:</p> <ol style="list-style-type: none"> 1. Validate that the implementation standard contains a data element which represents the same concept as that defined for the UML class. 2. Validate that the data element has the same relationships with other elements as those defined for the UML class. Validate that those relationships have the same source, target, direction, roles, and multiplicity as those documented in the Conceptual Model. 3. Validate that the data element has the same properties (attributes) as those specified for the UML class. Validate that those properties have the same name, definition, type, and multiplicity of those documented in the Conceptual Model. 4. Validate that the properties of the data element include those of all super classes of the UML class as documented in the Conceptual Model. Validate that those properties have the same name, definition, type, and multiplicity of those documented in the Conceptual Model. 5. Validate that the associations represented for the data element include those of all super classes of the UML class as documented in the Conceptual Model. Validate that those representations have the same source, target, roles, and multiplicity of those documented in the Conceptual Model. 6. Validate that the implementation standard enforces all constraints imposed on the UML class by the Conceptual Model.
-------------	--

A.3. Conformance class AI_TrainingData

Conformance class 2

<http://www.opengis.net/spec/TrainingDML-AI-1/1.0/conf/aitrainingdata>

Requirements class	/req/aitrainingdata
--------------------	---------------------

Abstract test A.2

/conf/aitrainingdata/classes

Requirement	/req/aitrainingdata/classes
Test purpose	To validate that the implementation standard correctly implements the UML Classes defined in the Conceptual Model.
Test method type	Manual Inspection

Abstract test A.2

Test method	<ol style="list-style-type: none">a. For each UML class defined or referenced in the AI_TrainingData Package:<ol style="list-style-type: none">1. Validate that the implementation standard contains a data element which represents the same concept as that defined for the UML class.2. Validate that the data element has the same relationships with other elements as those defined for the UML class. Validate that those relationships have the same source, target, direction, roles, and multiplicity as those documented in the Conceptual Model.3. Validate that the data element has the same properties (attributes) as those specified for the UML class. Validate that those properties have the same name, definition, type, and multiplicity of those documented in the Conceptual Model.4. Validate that the properties of the data element include those of all super classes of the UML class as documented in the Conceptual Model. Validate that those properties have the same name, definition, type, and multiplicity of those documented in the Conceptual Model.5. Validate that the associations represented for the data element include those of all super classes of the UML class as documented in the Conceptual Model. Validate that those representations have the same source, target, roles, and multiplicity of those documented in the Conceptual Model.6. Validate that the implementation standard enforces all constraints imposed on the UML class by the Conceptual Model.
-------------	--

A.4. Conformance class AI_Task

Conformance class 3

<http://www.opengis.net/spec/TrainingDML-AI-1/1.0/conf/aitask>

Requirements class /req/aitask

Abstract test A.3

/conf/aitask/classes

Requirement	/req/aitask/classes
Test purpose	To validate that the implementation standard correctly implements the UML Classes defined in the Conceptual Model.
Test method type	Manual Inspection

Abstract test A.3

Test method	<ol style="list-style-type: none">a. For each UML class defined or referenced in the AI_Task Package:<ol style="list-style-type: none">1. Validate that the implementation standard contains a data element which represents the same concept as that defined for the UML class.2. Validate that the data element has the same relationships with other elements as those defined for the UML class. Validate that those relationships have the same source, target, direction, roles, and multiplicity as those documented in the Conceptual Model.3. Validate that the data element has the same properties (attributes) as those specified for the UML class. Validate that those properties have the same name, definition, type, and multiplicity of those documented in the Conceptual Model.4. Validate that the properties of the data element include those of all super classes of the UML class as documented in the Conceptual Model. Validate that those properties have the same name, definition, type, and multiplicity of those documented in the Conceptual Model.5. Validate that the associations represented for the data element include those of all super classes of the UML class as documented in the Conceptual Model. Validate that those representations have the same source, target, roles, and multiplicity of those documented in the Conceptual Model.6. Validate that the implementation standard enforces all constraints imposed on the UML class by the Conceptual Model.
-------------	--

A.5. Conformance class AI_Label

Conformance class 4

<http://www.opengis.net/spec/TrainingDML-AI-1/1.0/conf/ailabel>

Requirements class	/req/ailabel
--------------------	--------------

Abstract test A.4

/conf/ailabel/classes

Requirement	/req/ailabel/classes
Test purpose	To validate that the implementation standard correctly implements the UML Classes defined in the Conceptual Model.
Test method type	Manual Inspection

Abstract test A.4

Test method	<ol style="list-style-type: none">a. For each UML class defined or referenced in the AI_Label Package:<ol style="list-style-type: none">1. Validate that the implementation standard contains a data element which represents the same concept as that defined for the UML class.2. Validate that the data element has the same relationships with other elements as those defined for the UML class. Validate that those relationships have the same source, target, direction, roles, and multiplicity as those documented in the Conceptual Model.3. Validate that the data element has the same properties (attributes) as those specified for the UML class. Validate that those properties have the same name, definition, type, and multiplicity of those documented in the Conceptual Model.4. Validate that the properties of the data element include those of all super classes of the UML class as documented in the Conceptual Model. Validate that those properties have the same name, definition, type, and multiplicity of those documented in the Conceptual Model.5. Validate that the associations represented for the data element include those of all super classes of the UML class as documented in the Conceptual Model. Validate that those representations have the same source, target, roles, and multiplicity of those documented in the Conceptual Model.6. Validate that the implementation standard enforces all constraints imposed on the UML class by the Conceptual Model.
-------------	---

A.6. Conformance class AI_Labeling

Conformance class 5

<http://www.opengis.net/spec/TrainingDML-AI-1/1.0/conf/ailabeling>

Requirements class /req/ailabeling

Abstract test A.5

/conf/ailabeling/classes

Requirement	/req/ailabeling/classes
Test purpose	To validate that the implementation standard correctly implements the UML Classes defined in the Conceptual Model.
Test method type	Manual Inspection

Abstract test A.5

Test method	<p>a. For each UML class defined or referenced in the AI_Labeling Package:</p> <ol style="list-style-type: none"> 1. Validate that the implementation standard contains a data element which represents the same concept as that defined for the UML class. 2. Validate that the data element has the same relationships with other elements as those defined for the UML class. Validate that those relationships have the same source, target, direction, roles, and multiplicity as those documented in the Conceptual Model. 3. Validate that the data element has the same properties (attributes) as those specified for the UML class. Validate that those properties have the same name, definition, type, and multiplicity of those documented in the Conceptual Model. 4. Validate that the properties of the data element include those of all super classes of the UML class as documented in the Conceptual Model. Validate that those properties have the same name, definition, type, and multiplicity of those documented in the Conceptual Model. 5. Validate that the associations represented for the data element include those of all super classes of the UML class as documented in the Conceptual Model. Validate that those representations have the same source, target, roles, and multiplicity of those documented in the Conceptual Model. 6. Validate that the implementation standard enforces all constraints imposed on the UML class by the Conceptual Model.
-------------	---

A.7. Conformance class AI_TDChangeset

Conformance class 6

<http://www.opengis.net/spec/TrainingDML-AI-1/1.0/conf/aitdchangeset>

Requirements class	/req/aitdchangeset
--------------------	--------------------

Abstract test A.6

/conf/aitdchangeset/classes

Requirement	/req/aitdchangeset/classes
Test purpose	To validate that the implementation standard correctly implements the UML Classes defined in the Conceptual Model.
Test method type	Manual Inspection

Abstract test A.6

Test method	<ol style="list-style-type: none">a. For each UML class defined or referenced in the AI_TDChangeset Package:<ol style="list-style-type: none">1. Validate that the implementation standard contains a data element which represents the same concept as that defined for the UML class.2. Validate that the data element has the same relationships with other elements as those defined for the UML class. Validate that those relationships have the same source, target, direction, roles, and multiplicity as those documented in the Conceptual Model.3. Validate that the data element has the same properties (attributes) as those specified for the UML class. Validate that those properties have the same name, definition, type, and multiplicity of those documented in the Conceptual Model.4. Validate that the properties of the data element include those of all super classes of the UML class as documented in the Conceptual Model. Validate that those properties have the same name, definition, type, and multiplicity of those documented in the Conceptual Model.5. Validate that the associations represented for the data element include those of all super classes of the UML class as documented in the Conceptual Model. Validate that those representations have the same source, target, roles, and multiplicity of those documented in the Conceptual Model.6. Validate that the implementation standard enforces all constraints imposed on the UML class by the Conceptual Model.
-------------	---

A.8. Conformance class AI_DataQuality

Conformance class 7

<http://www.opengis.net/spec/TrainingDML-AI-1/1.0/conf/aidataquality>

Requirements class /req/aidataquality

Abstract test A.7

/conf/aidataquality/classes

Requirement	/req/aidataquality/classes
Test purpose	To validate that the implementation standard correctly implements the UML Classes defined in the Conceptual Model.
Test method type	Manual Inspection

Abstract test A.7

Test method	<p>a. For each UML class defined or referenced in the AI_DataQuality Package:</p> <ol style="list-style-type: none">1. Validate that the implementation standard contains a data element which represents the same concept as that defined for the UML class.2. Validate that the data element has the same relationships with other elements as those defined for the UML class. Validate that those relationships have the same source, target, direction, roles, and multiplicity as those documented in the Conceptual Model.3. Validate that the data element has the same properties (attributes) as those specified for the UML class. Validate that those properties have the same name, definition, type, and multiplicity of those documented in the Conceptual Model.4. Validate that the properties of the data element include those of all super classes of the UML class as documented in the Conceptual Model. Validate that those properties have the same name, definition, type, and multiplicity of those documented in the Conceptual Model.5. Validate that the associations represented for the data element include those of all super classes of the UML class as documented in the Conceptual Model. Validate that those representations have the same source, target, roles, and multiplicity of those documented in the Conceptual Model.6. Validate that the implementation standard enforces all constraints imposed on the UML class by the Conceptual Model.
-------------	---

Appendix B: Example (Informative)

B.1. TrainingDataset encoding examples

B.1.1. WHU-RS19 dataset

The [WHU-RS19 dataset](#) is widely used in scene classification of remote sensing images. This dataset is collected from Google Earth and has 19 classes including airport, beach, bridge, commercial, desert, farmland, football field, forest, industrial, meadow, mountain, park, parking, pond, port, railway station, residential, river, and viaduct. Each class contains around 50 images, with the image size 600×600 and a resolution of 0.5m.

An example of JSON encoding of the WHU-RS19 dataset following the TrainingDML-AI UML model can be found in https://github.com/opengeospatial/TrainingDML-AI_SWG/tree/main/use-cases/examples/1.0/WHU-RS19.json.

B.1.2. DOTA-v1.5 dataset

The [DOTA-v1.5 dataset](#) is a large-scale dataset for object detection in aerial images. The sources for content in the dataset include Google Earth, Gaofen-2, and Jilin-1 imagery provided by China Resources Satellite Data Center. The 16 classes in DOTA-v1.5 are plane, ship, storage tank, baseball diamond, tennis court, basketball court, ground track field, harbor, bridge, large vehicle, small vehicle, helicopter, roundabout, soccer ball field, swimming pool, and container crane. Compared with other aerial image object detection datasets, the dataset has the largest number of classes. The images in the dataset have various image sizes (from 800×800 to 2000×2000) and resolutions (Google Earth/0.1m-1m, Gaofen-2/1m, Jilin-1/0.72m).

An example of JSON encoding of the DOTA-v1.5 dataset following the TrainingDML-AI UML model can be found in https://github.com/opengeospatial/TrainingDML-AI_SWG/tree/main/use-cases/examples/1.0/DOTA-v1.5.json.

B.1.3. KITTI 2D object detection dataset

The [KITTI 2D object detection dataset](#) is a novel open-access dataset and benchmark for road area and ego-lane detection. KITTI 2D consists of 7481 annotated training images of high variability from the KITTI autonomous driving platform by 2 PointGrey Flea2 color cameras, capturing a broad spectrum of urban street views and road scenes. The eight (8) classes in the KITTI 2D object detection dataset are car, van, truck, pedestrian, person_sitting, cyclist, tram, and misc. Compared with other street view object detection datasets, this dataset compresses diverse scenarios and captures real-world traffic situations, ranging from freeways over rural areas to inner-city scenes with many static and dynamic objects.

An example of JSON encoding of the KITTI 2D object detection dataset following the TrainingDML-AI UML model can be found in https://github.com/opengeospatial/TrainingDML-AI_SWG/tree/main/use-cases/examples/1.0/KITTI.json.

B.1.4. GID dataset

The [GID dataset](#) is one of start-of-art land cover classification datasets. This dataset has a large spatial coverage covering many provinces in China with a relatively high spatial resolution (2m). GID has two sets. One is the GID-5C. It has 150 images (image size 7200×6800) that are classified into 5 land cover classes. The other set is GID-15C. The images from GID-5C are sliced into 30,000 patches in GID-15C, which have three types of patch sizes (56×56, 112×112, 224×224) and are classified into 15 land cover classes.

An example of JSON encoding of the GID-5C dataset following the TrainingDML-AI UML model can be found in https://github.com/opengeospatial/TrainingDML-AI_SWG/tree/main/use-cases/examples/1.0/GID-5C.json.

B.1.5. Toronto3D dataset

The [Toronto3D dataset](#) is a large urban outdoor point cloud dataset for segmentation collected by the Mobile Laser Scanning System. The dataset covers about 1 km of scene streets in Toronto, including four areas named L001, L002, L003, and L004, with a total of 78.3 million points. Each point in this dataset has 10 attributes representing the 3D position, RGB color, intensity, GPS time, scan angle rank, and category, respectively. This dataset has eight categories, including road, road mark, natural, building, utility line, pole, car, and fence.

An example of JSON encoding of the Toronto3D dataset following the TrainingDML-AI UML model can be found in https://github.com/opengeospatial/TrainingDML-AI_SWG/tree/main/use-cases/examples/1.0/Toronto_3D.json.

B.1.6. WHU-Building dataset

The [WHU-Building dataset](#) is a change detection dataset collected from the Land Information New Zealand Data Service. The dataset is composed of images (with the resolution 0.2m) in 2012 and 2016, covering 20.5 km². It includes 12,796 and 16,077 buildings respectively in 2012 and 2016.

An example of JSON encoding of the WHU-Building dataset following the TrainingDML-AI UML model can be found in https://github.com/opengeospatial/TrainingDML-AI_SWG/tree/main/use-cases/examples/1.0/WHU-building.json.

B.1.7. California change detection dataset

The [California Change Detection Dataset](#) is composed of two images and a label image. The first image is a Landsat 8 acquisition covering Sacramento County, Yuba County and Sutter County, California, on 5 January 2017. It has nine channels covering the spectrum from deep blue to short-wave infrared, plus two long-wave infrared channels. The second image was acquired on 18 February 2017 by Sentinel-1A over the same area after the occurrence of a flood. The image is recorded in polarizations VV and VH and augmented with the ratio between the two intensities as a third channel. All these channels are log-transformed.

An example of JSON encoding of the California change detection dataset following the TrainingDML-AI UML model can be found in https://github.com/opengeospatial/TrainingDML-AI_SWG/tree/main/use-cases/examples/1.0/UiT_HCD_California_2017.json.

B.1.8. WHU MVS dataset

The WHU MVS dataset is a synthetic aerial dataset created for large-scale and high-resolution Earth surface reconstruction. The basic training sample of the dataset is a multi-view unit consisting of five aerial images, and their corresponding depth maps are taken as ground truth. There are a total of 5680 pairs of five-view aerial images in the dataset. All the images are simulated from a 3D surface model, which is produced by Smart3D software using Unmanned Aerial Vehicle (UAV) images and refined by manual editing.

An example of JSON encoding of the WHU MVS dataset following the TrainingDML-AI UML model can be found in https://github.com/opengeospatial/TrainingDML-AI_SWG/tree/main/use-cases/examples/1.0/WHU_MVS.json.

B.2. DataQuality encoding example

B.2.1. WHU-RS19 data quality

An encoded data quality example of the WHU-RS19 datasets following the TrainingDML-AI UML model can be found in https://github.com/opengeospatial/TrainingDML-AI_SWG/tree/main/use-cases/examples/1.0/WHU-RS19-quality.json.

B.3. TDChangeset encoding examples

B.3.1. DOTA-v1.5 changeset

DOTA-v1.5 uses the same images as DOTA-v1.0, but the extremely small instances (less than 10 pixels) are also annotated. Moreover, a new category “container crane” is added. It contains 403,318 instances in total. The number of images and dataset splits are the same as DOTA-v1.0. This version was released for the DOAI Challenge 2019 on Object Detection in Aerial Images in conjunction with IEEE CVPR 2019.

An encoded changeset example between the DOTA-v1.0 and DOTA-v1.5 datasets following the TrainingDML-AI UML model can be found in https://github.com/opengeospatial/TrainingDML-AI_SWG/tree/main/use-cases/examples/1.0/DOTA-v1.5-changeset.json.

Appendix C: Revision History (Informative)

Date	Release	Editor	Primary clauses modified	Description
2022-08-01	0.1	Peng Yue, Boyi Shangguan	all	Draft for internal review
2022-09-09	0.2	Peng Yue, Boyi Shangguan, Carl Reed	all	Merge edits and comments from Carl Reed.
2022-11-10	0.3	Peng Yue, Boyi Shangguan	all	Revision based on comments from Testbed-18 TrainingDataset.
2023-01-13	0.8	Peng Yue, Boyi Shangguan	all	Revision based on comments from Testbed-18 ML TD ER.
2023-01-29	1.0	Peng Yue, Boyi Shangguan, Ruixiang Liu	all	Internal review done and submission to request OAB review.

Appendix D: Bibliography

- [1] ISO, 2019. ISO 19107: 2019. Geographic information — Spatial schema. <https://www.iso.org/standard/26012.html>
- [2] ISO, 2022. ISO 19157-1: 2022. Geographic information — Data quality. <https://www.iso.org/standard/32575.html>
- [3] ISO, 2014. ISO 19115-1: 2014, Geographic information — Metadata — Part 1: Fundamentals. <https://www.iso.org/standard/53798.html>
- [4] ISO, 2015. ISO 19103: 2015, Geographic information — Conceptual schema language. <https://www.iso.org/standard/56734.html>
- [5] Landry, T., ed., 2018. OGC Testbed-14: Machine Learning Engineering Report, OGC 18-038r2. Wayland, MA: Open Geospatial Consortium Inc. <https://docs.ogc.org/per/18-038r2.html>
- [6] Meek, S., ed., 2019. OGC Testbed-15: Machine Learning Engineering Report, OGC 19-027r2. Wayland, MA: Open Geospatial Consortium Inc. <https://docs.ogc.org/per/19-027r2.html>
- [7] Schumann, G., ed., 2020. OGC Testbed-16: Machine Learning Training Data Engineering Report, OGC 20-018. Wayland, MA: Open Geospatial Consortium Inc. <https://docs.ogc.org/per/20-015r2.html>
- [8] Yue, P., Shangguan, B., Hu, L., Jiang, L., Zhang, C., Cao, Z., Pan, Y., 2022. Towards a training data model for artificial intelligence in earth observation. International Journal of Geographical Information Science, 1-25. <https://doi.org/10.1080/13658816.2022.2087223>