# Spatial data discovery and indexing tools: an approach based on metadata and fitness for use

**Pedro Castro**[1,3,4], Joaquim Alonso[2,3,4], Ivone Martins[3,4], João Honrado[3,5], Johannes Peterseil[6]

(1) Arc4DigiT - Applied Research Centre for Digital Transformation
(2) ProMetheus - Research Unit in Materials, Energy and Environment for Sustainability
(3) CIBIO-InBIO - Research Centre in Biodiversity and Genetic Resources
(4) Instituto Politécnico de Viana do Castelo
(5) Faculdade de Ciências, Universidade do Porto
(6) Umweltbundesamt Gmbh · Ecosystem Research and Environmental Information Management

# Overview

- Context

- Objectives

- Framework Overview

- Workflow

# Spatial Data Quality

Spatial Data Quality issues are important:

- the **increasing amount of spatial data production, handling and sharing** with different sources, different frequency of acquisition, different spatio-temporal scales, different levels of accuracy, different processing methods or techniques leads to many challenges in SDQ assessment

- Spatial Data is used **in very different application contexts** (data is often used with purposes other than producer' intended ones)

- It is necessary to **consider data quality** to identify datasets that satisfy the requirements of a particular application for specific user

Las IDE locales, acercando la información digital a los ciudadanos.

23, 24 y 25 de octubre

# Metadata

- Efforts have been made in metadata development and in meta-evaluation of external and (in)direct quality by the end-user(s), taking advantage of **metadata documentation possibilities and quality communication**

- Standard metadata profiles can contain a description of attributes about **dataset/database content, access and use conditions**, thus **allowing the assessment of data quality components and elements** (ISO 19157) as well as data quality management (ISO 19158).

- In this context, metadata catalogues offer opportunities for the implementation and improvement of spatial data quality evaluation/assessment tools related to knowledge discovery, searching and indexing

# ThemisE

THEmatic Metadata-based and fItness-for-use Spatial data quality Evaluation platform – **ThemisE** platform:

- implemented as an autonomous and modular **Web application** to perform **quality evaluation BASED ON METADATA** considering that
  - Metadata can contain information about the content, quality, condition and other characteristics of the data (ISO 2005) that can be used for (meta)quality evaluation
  - Frequent limitations to data access and use
  - Increasing availability of metadata catalogues allowing a (simple) integration with an evaluation platform

- With the aim to support the **quality-driven discovery and selection of relevant datasets** (or the identification of data gaps) necessary for environmental/ecological modelling based on (well documented) datasets' metadata

Las IDE locales, acercando la información digital a los ciudadanos.

23, 24 y 25 de octubre

# Framework overview

ThemisE platform allows two types of evaluation:

- an internal evaluation centered on the comparison of the characteristics of the dataset, as detailed in metadata by the producer, with the required elements according to a predefined standard profile;

- an external quality evaluation that is based on determining the **matching level** (fitness-for-use) between the **characteristics of the dataset** (detailed by its **metadata**) and the **characteristics of the data required by the user** that describe the user's requirements for a given application context (and defined through expected values for predefined quality indicators)

focused on evaluating how data will fit the users' needs, to bring data sets closer to users' applications

Las IDE locales, acercando la información digital a los ciudadanos.

23, 24 y 25 de octubre

# General process



2. Extract internal data quality values from metadata catalogues for quality indicators selected by the user

3. Run comparison of information provided by users and metadata

Metadata catalogue(s) (e.g. catalogue from task 5.7)

DEFINE

Configuration

Data on quality indicators and respective evaluation rules

Data on metadata/quality elements

ACCESS

EVALUATION PARAMETERS

User-oriented data quality evaluation routines

CONTROLS

USER

VALUES FOR EVALUATION

PRODUCE

Expected quality indicators values (user's requirements) for a specific application context

DEFINE

Report(s) on data quality

1. Users specify expected data quality values for a set of pre-defined quality indicators relevant to the application context

# Functional workflow

# Actions workflow

# Thematic category



**Add Thematic Category**

Allows to define different datasets that are targeted for search for the given application context

| Quality Indicator | Land use ✖ | Orthoimagery ✖ | |
|---|---|---|---|
| Description | TC's description | TC's description | |
| Filter by abstract/title | 1 value ≡ | 1 value ≡ | |
| Topic category | 1 value ≡ | 1 value ≡ | |
| Spatial scale | 1 value ≡ | 1 value ≡ | |
| Spatial extent | 1 value | ✚ | |

ECOPOTENTIAL Project    inBIO    ThemisE - THEmatic Metadata-based and fitness for use Spatial data quality Evaluation    No 641762

**Show evaluation's result**

# Quality indicators



Users specify expected data quality values of available quality indicators for each targeted dataset

# Metadata catalogues

# Evaluation

# Results

Las IDE locales, acercando la información digital a los ciudadanos.

23, 24 y 25 de octubre

# Thank you for your attention