# ESCAPE

European Science Cluster of Astronomy & Particle physics ESFRI research Infrastructures

# How to harvest and retrieve entries from the OSSR

Enrique GARCIA & Thomas VUILLAUME
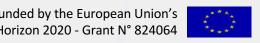
FG2 call – 26/05/2021

# Outline

- Harvest metadata from a repository
  - OAI-PMH protocol
  - Repository metadata representations Vs. entry metadata

- Query entry metadata
  - Zenodo REST API
  - OAI-PMH Vs. REST API

- Demo

# Harvesting metadata from a repository

- The provider is responsible of the service.
  - It chooses the harvest protocol and metadata representation of the records.

- All metadata records are licenced under CC0 license, besides the license applying to the data files of the digital entry.

- Zenodo:
  - OAI-PMH protocol
    - https://developers.zenodo.org/#oai-pmh
  - Metadata supported schemas:
    - DataCite (various versions)
    - Dublin Core
    - MARC21

# OAI-PMH Protocol

- Open Archives Initiative Protocol for Metadata Harvesting
  - Developed for harvesting records in an archive/repository.

- Uses a base URL to which you can add different *"verbs"* to reduce the query/search.
  - OAI-PMH (exhaustive!) tutorial
    https://indico.cern.ch/event/5710/sessions/108048/
    attachments/988151/1405129/Simeon_tutorial.pdf

- Various 'harvester'-libraries in various programming languages
  - Python: oai-harvest, pyoai…

## Six verbs

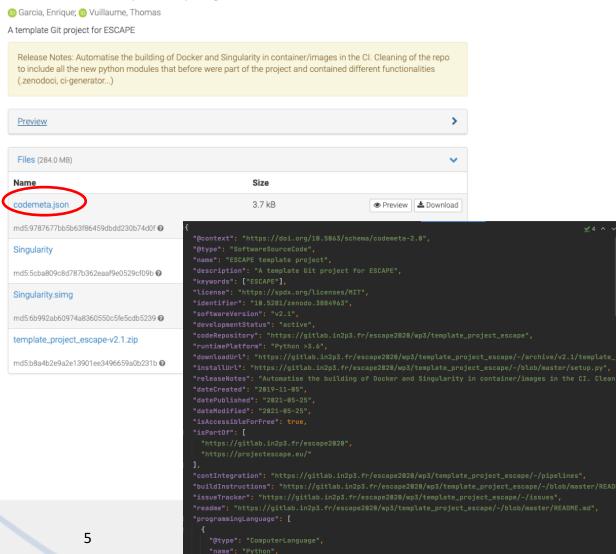| | Verb | Function |
|---|---|---|
| metadata about the repository | Identify | Description of repository |
| | ListMetadataFormats | Metadata formats supported by repository |
| | ListSets | Sets defined by repository |
| harvesting verbs | ListIdentifiers | List OAI unique identifiers contained in repository |
| | ListRecords | List of many records |
| | GetRecord | List a single record |

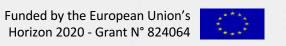# Repository metadata representation Vs. resource metadata

🟡 Not the same !

# Query metadata within the same record

- **Zenodo REST API**
  - allows accessing to this metadata, if available.
  - License ?

- **OAI-PMH Vs. REST API**
  - AOI – PMH
    - **+** Better harvesting
      - Faster.
      - Thought for large and continuous queries.
    - **−** Metadata representation provided by the data provider.
  - REST API
    - **+** Access to the full record information.
    - **−** Harvesting not optimized for large searches.

# Demo

## REST API handlers (python)
- ZenodoCI: https://gitlab.in2p3.fr/escape2020/wp3/zenodoci
- PyZenodo: https://github.com/space-physics/pyzenodo3
- zenodo-python: https://github.com/SiLeBAT/zenodo-python

## OAI-PMH harvesters (python)
- OAI-PMH Harvest: https://github.com/bloomonkey/oai-harvest
- pyoai: https://github.com/infrae/pyoai