
EURITO Documentation

eurito

Nov 28, 2019

Contents:

1	Data sources and core processing	3
1.1	CORDIS	3
1.2	arXiv	6
1.3	PATSTAT	6
1.4	Companies	6
2	Batchables	9
2.1	run.py (lolvelty)	9
3	Production pipelines	11
3.1	Transfer of Elasticsearch data	11
3.2	Centrality Pipeline	12
3.3	Cordis to Neo4j	13
3.4	OpenAIRE to Neo4j	14
4	Ontologies and schemas	15
4.1	Tier 0	15
4.2	Tier 1	15
5	Scripts	17
5.1	nesta_prepare_batch	17
5.2	nesta_docker_build	17
6	License	19
	Python Module Index	21
	Index	23

Branch	Docs	Build
Master		
Development		

Welcome to **EURITO**! This repository houses our fully-audited tools and packages, as well as our in-house production system. If you're reading this on [our GitHub repo](#), you will find complete documentation at our [Read the Docs site](#).

CHAPTER 1

Data sources and core processing

In EURITO we predominantly use four data sources:

- EU-funded research from *CORDIS*
- Technical EU research on *arXiv*
- EU Patents from *PATSTAT*
- EU Companies [under license, we can't specify the source publicly. Contact us for more details!]

1.1 CORDIS

Data from the CORDIS's [H2020 API](#) and [FP7 API](#) funded projects is extracted using code found [in this repository](#).

In total, 51250 organisations and 50640 projects were extracted from the API. There are 1102 proposal calls, 245465 publications and 34507 reports. In total 6545 are associated with the projects.

Software outputs are associated with the projects, using [OpenAIRE API](#).

All of these entities are then linked together, and stored using a [neo4j](#) graph database. The code for automatically piping the data in neo4j is provided [here](#).

1.1.1 Cordis to Neo4j

Tools for piping data from a SQLAlchemy ORM to Neo4j, to be used in the Luigi pipeline.

`orm_to_neo4j(session, transaction, orm_instance, parent_orm=None, rel_name=None)`

Pipe a SQLAlchemy ORM instance (a ‘row’ of data) to neo4j, inserting it as a node or relationship, as appropriate.

Parameters

- `session(sqlalchemy.Session)` – SQL DB session.
- `transaction(py2neo.Transaction)` – Neo4j transaction

- **orm_instance** (*sqlalchemy.Base*) – Instance of a SQLAlchemy ORM, i.e. a ‘row’ of data.
- **parent_orm** (*sqlalchemy.Base*) – Parent ORM to build relationship to
- **rel_name** (*str*) – Name of the relationship to be added to Neo4j

build_relationships (*session, graph, orm, data_row, rel_name, parent_orm=None*)

Build a py2neo.Relationship object from SQLAlchemy objects.x

Parameters

- **session** (*sqlalchemy.Session*) – SQL DB session.
- **transaction** (*py2neo.Transaction*) – Neo4j transaction
- **orm** (*sqlalchemy.Base*) – A SQLAlchemy ORM
- **rel_name** (*str*) – Name of the relationship to be added to Neo4j
- **parent_orm** (*sqlalchemy.Base*) – Another ORM to build relationship to. If this is not specified, it implies that **orm** is node, rather than a relationship.

Returns Relationships pointing to the node (inferred from ORM), and one pointing back to its associated project.

Return type {relationship, back_relationship}

set_constraints (*orm, graph_schema*)

Set constraints in the neo4j graph schema.

Parameters

- **orm** (*sqlalchemy.Base*) – A SQLAlchemy ORM
- **graph_schema** (*py2neo.Graph.Schema*) – Neo4j graph schema.

prepare_base_entities (*table*)

Returns the objects required to generate a graph representation of the ORM.

Parameters **table** (*sqlalchemy.sql.Table*) – SQL alchemy table object from which to extract an graph representation.

Returns

Two ORMs and a string describing their relationship

Return type {orm, parent_orm, rel_name}

flatten (*orm_instance*)

Convert a SQLAlchemy ORM (i.e. a ‘row’ of data) to flat JSON.

Parameters **orm_instance** (*sqlalchemy.Base*) – Instance of a SQLAlchemy ORM, i.e. a ‘row’ of data.

Returns A flat row of data, inferred from *orm_instance*

Return type row (dict)

flatten_dict (*row, keys=[('title',), ('street', 'city', 'postalCode')]*)

Flatten a dict by concatenating string values of matching keys.

Parameters **row** (*dict*) – Data to be flattened

Returns Concatenated data.

Return type flat (str)

retrieve_node (*session, graph, orm, parent_orm, data_row*)

Retrieve an existing node from neo4j, by first retrieving it's id (field name AND value) via SQLAlchemy.

Parameters

- **session** (*sqlalchemy.Session*) – SQL DB session.
- **transaction** (*py2neo.Transaction*) – Neo4j transaction
- **orm** (*sqlalchemy.Base*) – SQLAlchemy ORM describing *data_row*
- **parent_orm** (*sqlalchemy.Base*) – Parent ORM to build relationship to
- **data_row** (*dict*) – Flat row of data retrieved from *orm*

Returns Node of data corresponding to *data_row*

Return type node (*py2neo.Node*)

table_from_fk (*fks*)

Get the table name of the fk constraint, ignoring the cordis_projects table

Parameters **fks** (list of *SqlAlchemy.ForeignKey*) – All foreign keys for a given table.

Returns The table name corresponding to the non-Project foreign key.

Return type tablename (str)

get_row (*session, parent_orm, orm, data_row*)

Retrieve a flat row of data corresponding to the parent relation, inferred via foreign keys.

Parameters

- **session** (*sqlalchemy.Session*) – SQL DB session.
- **parent_orm** (*sqlalchemy.Base*) – Parent ORM to build relationship to
- **orm** (*sqlalchemy.Base*) – SQLAlchemy ORM describing *data_row*
- **data_row** (*dict*) – Flat row of data retrieved from *orm*

Returns Flat row of data retrieved from *parent_orm*

Return type _row (dict)

1.1.2 Enrich Cordis with OpenAIRE

Tools for collecting OpenAIRE data (by Cordis project), and piping to Neo4j.

write_record_to_neo (*record, output_type, graph*)

A utility function, which takes record and writes it to neo4j graph

Parameters

- **record** (*dict*) – a dictionary that contains metadata about a record
- **output_type** (str) – type of record to be extracted from OpenAIRE API. Accepts “software”, “datasets”, “publications”, “ECProjects”
- **graph** (*graph_session*) – connection to neo4j database

get_project_soups (*currentUrl, reqsession, output_type, projectID*)

Gets a beautiful soup according to output type and projectID

Parameters

- **currentUrl** (str) – URL to OpenAIRE API

- **reqsession** (*instance of Requests session*) – currently open HTTP request
- **output_type** (*str*) – type of record to be extracted from OpenAIRE API. Accepts “software”, “datasets”, “publications”, “ECProjects”
- **projectID** (*str*) – EC project identifier

Returns a list of BeautifulSoup objects that contain the results from API call

Return type souplist(list)

get_results_from_soups (*souplist*)

Extracts string from all BeautifulSoup objects and merges them into one list

Parameters **souplist** (*list*) – a list of BeautifulSoup objects that contain the results from API call

Returns a list of strings with results metadata

Return type resultlist(list)

1.2 arXiv

All articles from [arXiv](#), which is the world’s premier repository of pre-prints of articles in the physical, quantitative and computational sciences, are already automatically collected, geocoded (using [GRID](#)) and enriched with topics (using [MAG](#)). Articles are assigned to EU NUTS regions (at all levels) using [nesta’s nuts-finder python package](#).

Data is transferred to EURITO’s [elasticsearch](#) server via [nesta’s es2es package](#). The [lolvley algorithm](#) is then applied to the data in order to generate a novelty metric for each article. This procedure is better described in [this blog](#) (see “Defining novelty”).

The indicators using this data source are presented in [this other EURITO repository](#).

In total, 1598033 articles have been processed, of which 459371 have authors based in EU nations.

1.3 PATSTAT

All patents from the PATSTAT service have been collected in [nesta’s own database](#) using [nesta’s pypatstat library](#). Since this database is very large, we have selected patents which belong to a patent family with a granted patent first published after the year 2000, with at least one person or organisation (inventor or applicant) based in an EU member state. This leads to 1552303 patents in the database.

Data is transferred to EURITO’s [elasticsearch](#) server via [nesta’s es2es package](#). The [lolvley algorithm](#) is then applied to the data in order to generate a novelty metric for each article. This procedure is better described in [this blog](#) (see “Defining novelty”).

The indicators using this data source are presented in [this other EURITO repository](#).

1.4 Companies

We have acquired private-sector company data under license. The dataset contains 550540 companies, of which 133641 are based in the EU.

Data is transferred to EURITO’s [elasticsearch](#) server via [nesta’s es2es package](#). The [lolvley algorithm](#) is then applied to the data in order to generate a novelty metric for each article. This procedure is better described in [this blog](#) (see “Defining novelty”).

The indicators using this data source are presented in [this other EURITO repository](#).

CHAPTER 2

Batchables

2.1 run.py (lolvelty)

Calculates the “lolvelty” novelty score to documents in Elasticsearch, on a document-by-document basis. Note that this is a slow procedure, and the bounds of document “lolvelty” can’t be known a priori.

`run()`

CHAPTER 3

Production pipelines

We use luigi routines to orchestrate our pipelines. The batching procedure relies on batchables as described in [batchables](#). Other than [luighacks.autobatch](#), which is respectively documented in Nesta's codebase, the routine procedure follows the [Luigi](#) documentation well.

3.1 Transfer of Elasticsearch data

This pipeline is responsible for the transfer of Elasticsearch data from a remote origin (in our case, Nesta's Elasticsearch endpoint) to EURITO's endpoint.

```
class Es2EsTask(*args, **kwargs)
    Bases: luigi.task.Task

    date = <luigi.parameter.DateParameter object>
    origin_endpoint = <luigi.parameter.Parameter object>
    origin_index = <luigi.parameter.Parameter object>
    dest_endpoint = <luigi.parameter.Parameter object>
    dest_index = <luigi.parameter.Parameter object>
    test = <luigi.parameter.BoolParameter object>
    chunksize = <luigi.parameter.IntParameter object>
    do_transfer_index = <luigi.parameter.BoolParameter object>
    db_config_path = <luigi.parameter.Parameter object>

    output()
        Points to the output database engine

    run()
        The task run method, to be overridden in a subclass.

        See Task.run
```

```
class EsLoveltyTask(*args, **kwargs)
Bases: nesta.core.luigihacks.estask.LazyElasticsearchTask

date = <luigi.parameter.DateParameter object>
origin_endpoint = <luigi.parameter.Parameter object>
origin_index = <luigi.parameter.Parameter object>
test = <luigi.parameter.BoolParameter object>
process_batch_size = <luigi.parameter.IntParameter object>
do_transfer_index = <luigi.parameter.BoolParameter object>

requires()
The Tasks that this Task depends on.
```

A Task will only run if all of the Tasks that it requires are completed. If your Task does not require any other Tasks, then you don't need to override this method. Otherwise, a subclass can override this method to return a single Task, a list of Task instances, or a dict whose values are Task instances.

See Task.requires

```
class RootTask(*args, **kwargs)
Bases: luigi.task.WrapperTask

production = <luigi.parameter.BoolParameter object>
date = <luigi.parameter.DateParameter object>

requires()
The Tasks that this Task depends on.
```

A Task will only run if all of the Tasks that it requires are completed. If your Task does not require any other Tasks, then you don't need to override this method. Otherwise, a subclass can override this method to return a single Task, a list of Task instances, or a dict whose values are Task instances.

See Task.requires

3.2 Centrality Pipeline

Takes network from Neo4j database, calculates network centrality measures and updates each node in the database with new centrality attributes

```
class RootTask(*args, **kwargs)
Bases: luigi.task.WrapperTask

The root task, which collects the supplied parameters and calls the main task.
```

Parameters

- **date** (*datetime*) – Date used to label the outputs
- **output_type** (*str*) – type of record to be extracted from OpenAIRE API. Accepts “software”, “datasets”, “publications”, “ECProjects”
- **production** (*bool*) – test mode or production mode

```
date = <luigi.parameter.DateParameter object>
output_type = <luigi.parameter.Parameter object>
production = <luigi.parameter.BoolParameter object>
```

```
requires()
    Call the task to run before this in the pipeline.

class CalcCentralityTask(*args, **kwargs)
    Bases: luigi.task.Task

    Takes network from Neo4j database, calculates network centrality measures and updates each node in the
    database with new centrality attributes

Parameters
    • date (datetime) – Date used to label the outputs
    • output_type (str) – type of record to be extracted from OpenAIRE API. Accepts “soft-
        ware”, “datasets”, “publications”, “ECProjects”
    • test (bool) – run a shorter version of the task if in test mode

date = <luigi.parameter.DateParameter object>
output_type = <luigi.parameter.Parameter object>
test = <luigi.parameter.BoolParameter object>

output()
    Points to the output database engine where the task is marked as done. The luigi_table_updates table exists
    in test and production databases.

run()
    The task run method, to be overridden in a subclass.

    See Task.run
```

3.3 Cordis to Neo4j

Task for piping Cordis data from SQL to Neo4j.

```
class CordisNeo4jTask(*args, **kwargs)
    Bases: luigi.task.Task

    Task for piping Cordis data to neo4j

    test = <luigi.parameter.BoolParameter object>
    date = <luigi.parameter.DateParameter object>

    output()
        Points to the output database engine where the task is marked as done. The luigi_table_updates table exists
        in test and production databases.

    run()
        The task run method, to be overridden in a subclass.

        See Task.run

class RootTask(*args, **kwargs)
    Bases: luigi.task.WrapperTask

    production = <luigi.parameter.BoolParameter object>

requires()
    The Tasks that this Task depends on.
```

A Task will only run if all of the Tasks that it requires are completed. If your Task does not require any other Tasks, then you don't need to override this method. Otherwise, a subclass can override this method to return a single Task, a list of Task instances, or a dict whose values are Task instances.

See Task.requires

3.4 OpenAIRE to Neo4j

Pipe data directly from the OpenAIRE API to Neo4j by matching to Cordis projects already in Neo4j.

```
class RootTask(*args, **kwargs)
Bases: luigi.task.WrapperTask
```

The root task, which collects the supplied parameters and calls the SimpleTask.

Parameters

- **date** (`datetime`) – Date used to label the outputs
- **output_type** (`str`) – type of record to be extracted from OpenAIRE API. Accepts “software”, “datasets”, “publications”, “ECProjects”
- **production** (`bool`) – test mode or production mode

```
date = <luigi.parameter.DateParameter object>
output_type = <luigi.parameter.Parameter object>
production = <luigi.parameter.BoolParameter object>
```

requires()

Call the task to run before this in the pipeline.

```
class OpenAireToNeo4jTask(*args, **kwargs)
Bases: luigi.task.Task
```

Takes OpenAIRE entities from MySQL database and writes them into Neo4j database

Parameters

- **date** (`datetime`) – Date used to label the outputs
- **output_type** (`str`) – type of record to be extracted from OpenAIRE API. Accepts “software”, “datasets”, “publications”, “ECProjects”
- **test** (`bool`) – run a shorter version of the task if in test mode

```
date = <luigi.parameter.DateParameter object>
output_type = <luigi.parameter.Parameter object>
test = <luigi.parameter.BoolParameter object>
```

output()

Points to the output database engine where the task is marked as done. The `luigi_table_updates` table exists in test and production databases.

run()

The task run method, to be overridden in a subclass.

See Task.run

CHAPTER 4

Ontologies and schemas

4.1 Tier 0

Raw data collections (“tier 0”) in the production system do not adhere to a fixed schema or ontology, but instead have a schema which is very close to the raw data. Modifications to field names tend to be quite basic, such as lowercase and removal of whitespace in favour of a single underscore.

4.2 Tier 1

Processed data (“tier 1”) is intended for public consumption, using a common ontology. The convention we use is as follows:

- Field names are composed of up to three terms: a `firstName`, `middleName` and `lastName`
- Each term (e.g. `firstName`) is written in lowerCamelCase.
- `firstName` terms correspond to a restricted set of basic quantities.
- `middleName` terms correspond to a restricted set of modifiers (e.g. adjectives) which add nuance to the `firstName` term. Note, the special `middleName` term of is reserved as the default value in case no `middleName` is specified.
- `lastName` terms correspond to a restricted set of entity types.

Valid examples are `date_start_project` and `title_of_project`.

Tier 0 fields are implicitly excluded from tier 1 if they are missing from the `schema_transformation` file. Tier 1 schema field names are applied via `nesta.packages.decorator.schema_transform`

CHAPTER 5

Scripts

A set of helper scripts for the batching system.

Note that this directory is required to sit in `$PATH`. By convention, all executables in this directory start with `nesta_` so that our developers know where to find them.

5.1 `nesta_prepare_batch`

Collect a batchable `run.py` file, including dependencies and an automatically generated requirements file; which is all zipped up and sent to AWS S3 for batching. This script is executed automatically in `luigihacks.autobatch.AutoBatchTask.run`.

Parameters:

- **BATCHABLE_DIRECTORY**: The path to the directory containing the batchable `run.py` file.
- **ARGS**: Space-separated-list of files or directories to include in the zip file, for example imports.

5.2 `nesta_docker_build`

Build a docker environment and register it with the AWS ECS container repository.

Parameters:

- **DOCKER_RECIPE**: A docker recipe. See `docker_recipes/` for a good idea of how to build a new environment.

CHAPTER 6

License

MIT © 2019 EURITO

Python Module Index

e

```
eurito_daps.batchables.novelty.lolvelty.run,  
    9  
eurito_daps.packages.cordis.cordis_neo4j,  
    3  
eurito_daps.packages.utils.openaire_utils,  
    5  
eurito_daps.routines.centrality.centrality,  
    12  
eurito_daps.routines.cordis.cordis_neo4j_task,  
    13  
eurito_daps.routines.es_data.es_data,  
    11  
eurito_daps.routines.openaire.openaire_to_neo4j_search,  
    14
```

Index

B

build_relationships() (in module `eurito_daps.packages.cordis.cordis_neo4j`), 4

C

`CalcCentralityTask` (class in `eurito_daps.routines.centrality.centrality`), 13
`chunksize` (`Es2EsTask` attribute), 11
`CordisNeo4jTask` (class in `eurito_daps.routines.cordis.cordis_neo4j_task`), 13

`date` (`CalcCentralityTask` attribute), 13
`date` (`CordisNeo4jTask` attribute), 13
`date` (`Es2EsTask` attribute), 11
`date` (`EsLoveltyTask` attribute), 12
`date` (`OpenAireToNeo4jTask` attribute), 14
`date` (`RootTask` attribute), 12, 14
`db_config_path` (`Es2EsTask` attribute), 11
`dest_endpoint` (`Es2EsTask` attribute), 11
`dest_index` (`Es2EsTask` attribute), 11
`do_transfer_index` (`Es2EsTask` attribute), 11
`do_transfer_index` (`EsLoveltyTask` attribute), 12

E

`Es2EsTask` (class in `eurito_daps.routines.es_data.es_data`), 11
`EsLoveltyTask` (class in `eurito_daps.routines.es_data.es_data`), 11
`eurito_daps.batchables.novelty.lovelty`.
 (`module`), 9
`eurito_daps.packages.cordis.cordis_neo4j`.
 (`module`), 3
`eurito_daps.packages.utils.openaire_utils`.
 (`module`), 5
`eurito_daps.routines.centrality.centrality`.
 (`module`), 12

`eurito_daps.routines.cordis.cordis_neo4j_task`.
 (`module`), 13
`eurito_daps.routines.es_data.es_data`.
 (`module`), 11
`eurito_daps.routines.openaire.openaire_to_neo4j_sea`.
 (`module`), 14

F

`flatten()` (in module `eurito_daps.packages.cordis.cordis_neo4j`), 4
`flatten_dict()` (in module `eurito_daps.packages.cordis.cordis_neo4j`), 4

G

`get_project_soups()` (in module `eurito_daps.packages.utils.openaire_utils`), 5
`get_results_from_soups()` (in module `eurito_daps.packages.utils.openaire_utils`), 6
`get_row()` (in module `eurito_daps.packages.cordis.cordis_neo4j`), 5

O

`OpenAireToNeo4jTask` (class in `eurito_daps.routines.openaire.openaire_to_neo4j_search`), 14
`origin_endpoint` (`Es2EsTask` attribute), 11
`origin_endpoint` (`EsLoveltyTask` attribute), 12
`origin_index` (`Es2EsTask` attribute), 11
`origin_index` (`EsLoveltyTask` attribute), 12
`orm_to_neo4j()` (in module `eurito_daps.packages.cordis.cordis_neo4j`), 3
`output()` (`CalcCentralityTask` method), 13
`output()` (`CordisNeo4jTask` method), 13
`output()` (`Es2EsTask` method), 11
`output()` (`OpenAireToNeo4jTask` method), 14

output_type (*CalcCentralityTask attribute*), 13
output_type (*OpenAireToNeo4jTask attribute*), 14
output_type (*RootTask attribute*), 12, 14

P

prepare_base_entities() (in module eu-
rito_daps.packages.cordis.cordis_neo4j),
4
process_batch_size (*EsLoveltyTask attribute*), 12
production (*RootTask attribute*), 12–14

R

requires () (*EsLoveltyTask method*), 12
requires () (*RootTask method*), 12–14
retrieve_node() (in module eu-
rito_daps.packages.cordis.cordis_neo4j),
4
RootTask (class in eu-
rito_daps.routines.centrality.centrality), 12
RootTask (class in eu-
rito_daps.routines.cordis.cordis_neo4j_task),
13
RootTask (class in eu-
rito_daps.routines.es_data.es_data), 12
RootTask (class in eu-
rito_daps.routines.openaire.openaire_to_neo4j_search),
14
run () (*CalcCentralityTask method*), 13
run () (*CordisNeo4jTask method*), 13
run () (*Es2EsTask method*), 11
run() (in module eu-
rito_daps.batchables.novelty.lovelty.run),
9
run () (*OpenAireToNeo4jTask method*), 14

S

set_constraints() (in module eu-
rito_daps.packages.cordis.cordis_neo4j),
4

T

table_from_fk() (in module eu-
rito_daps.packages.cordis.cordis_neo4j),
5
test (*CalcCentralityTask attribute*), 13
test (*CordisNeo4jTask attribute*), 13
test (*Es2EsTask attribute*), 11
test (*EsLoveltyTask attribute*), 12
test (*OpenAireToNeo4jTask attribute*), 14

W

write_record_to_neo() (in module eu-
rito_daps.packages.utils.openaire_utils),
5