

---

# nesta Documentation

nesta

May 07, 2020



---

## Contents:

---

<b>1 Packages</b>	<b>3</b>
1.1 Code and scripts . . . . .	3
1.1.1 Meetup . . . . .	3
1.1.2 Health data . . . . .	7
1.1.3 NLP Utils . . . . .	8
1.1.4 Geo Utils . . . . .	8
1.1.5 Format Utils . . . . .	9
1.1.6 Decorators . . . . .	10
1.2 Code auditing . . . . .	11
<b>2 Production</b>	<b>13</b>
2.1 How to put code into production at nesta . . . . .	13
2.2 Code and scripts . . . . .	14
2.2.1 Routines . . . . .	14
2.2.2 Batchables . . . . .	46
2.2.3 ORMs . . . . .	52
2.2.4 Ontologies and schemas . . . . .	55
2.2.5 Luigi Hacks . . . . .	56
2.2.6 Scripts . . . . .	61
2.2.7 Scripts . . . . .	61
2.2.8 Elasticsearch . . . . .	62
2.2.9 Containerised Luigi . . . . .	63
<b>3 AWS FAQ</b>	<b>67</b>
3.1 Where is the data? . . . . .	67
3.2 Why don't you use Aurora rather than MySQL? . . . . .	67
3.3 Where are the production machines? . . . . .	67
3.4 Where is the latest config? . . . . .	67
3.5 Where do I start with Elasticsearch? . . . . .	68
<b>4 Troubleshooting</b>	<b>69</b>
4.1 How do I restart the apache server after downtime? . . . . .	69
4.2 How do I restart the luigi server after downtime? . . . . .	69
4.3 How do I perform initial setup to ensure the batchables will run? . . . . .	69
4.4 How do I add a new user to the server? . . . . .	69
<b>5 Packages</b>	<b>71</b>

---

<b>6 Production</b>	<b>73</b>
6.1 License . . . . .	73
<b>Python Module Index</b>	<b>75</b>
<b>Index</b>	<b>77</b>

Branch	Docs	Build
Master		
Development		

Welcome to [nesta!](#) 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

---

## Packages

---

Nesta's collection of tools for meaty tasks. Any processes that go into production come here first, but there are other good reasons for code to end up here.

### 1.1 Code and scripts

#### 1.1.1 Meetup

**NB: The meetup pipeline will not work until this issue has been resolved.**

Data collection of Meetup data. The procedure starts with a single country and [Meetup category](#). All of the groups within the country are discovered, from which all members are subsequently retrieved (no personal information!). In order to build a fuller picture, all other groups to which the members belong are retrieved, which may be in other categories or countries. Finally, all group details are retrieved.

The code should be executed in the following order, which reflects the latter procedure:

- 1) country\_groups.py
- 2) groups\_members.py
- 3) members\_groups.py
- 4) groups\_details.py

Each script generates a list of dictionaries which can be ingested by the proceeding script.

#### Country → Groups

Start with a country (and Meetup category) and end up with Meetup groups.

**generate\_coords** ( $x0, y0, x1, y1, n$ )

Generate  $\mathcal{O}(\frac{n^2}{2})$  coordinates in the bounding box  $(x0, y0), (x1, y1)$ , such that overlapping circles of equal radii (situated at each coordinate) entirely cover the area of the bounding box. The longitude and latitude are treated

as euclidean variables, although the radius (calculated from the smallest side of the bounding box divided by  $n$ ) is calculated correctly. In order for the circles to fully cover the region, an unjustified factor of 10% is included in the radius. Feel free to do the maths and work out a better strategy for covering a geographical area with circles.

The circles (centred on each X) are staggered as so (single vertical lines or four underscores correspond to a circle radius):

```
____X____ ____X____  
|  
X_____X_____X  
|  
____X____ ____X____
```

This configuration corresponds to  $n = 4$ .

### Parameters

- **x0, y0, x1, y1** (*float*) – Bounding box coordinates (lat/lon)
- **n** (*int*) – The fraction by which to calculate the Meetup API radius parameter, with respect to the smallest side of the country's shape bbox. This will generate  $\mathcal{O}(\frac{n^2}{2})$  separate Meetup API radius searches. The total number of searches scales with the ratio of the bbox sides.

**Returns** The radius and coordinates for the Meetup API request

**Return type** float, list of tuple

### get\_coordinate\_data (n)

Generate the radius and coordinate data (see `generate_coords`) for each shape (country) in the shapefile pointed to by the environmental variable `WORLD_BORDERS`.

**Parameters** **n** (*int*) – The fraction by which to calculate the Meetup API radius parameter, with respect to the smallest side of the country's shape bbox. This will generate  $\mathcal{O}(\frac{n^2}{2})$  separate Meetup API radius searches. The total number of searches scales with the ratio of the bbox sides.

**Returns**

**containing coordinate and radius** for each country.

**Return type** pd.DataFrame

### assert\_iso2\_key (df, iso2)

### class MeetupCountryGroups (country\_code, coords, radius, category, n=10)

Bases: object

Extract all meetup groups for a given country.

#### country\_code

ISO2 code

**Type** str

#### params (

obj:'dict'): GET request parameters, including lat/lon.

#### groups

List of meetup groups in this country, assigned after calling `get_groups`.

**Type** list of str

**get\_groups** (*lon, lat, offset=0, max\_pages=None*)

Recursively get all groups for the given parameters. It is assumed that you will run with the default arguments, since they are set automatically in the recursing procedure.

**get\_groups\_recursive()**

Call `get_groups` for each lat,lon coordinate

**Groups → Members**

Start with Meetup groups and end up with Meetup members.

**get\_members** (*params*)

Hit the Meetup API for the members of a specified group.

**Parameters** **params** (dict) – <https://api.meetup.com/members/> parameters

**Returns** Meetup member IDs

**Return type** (list of str)

**get\_all\_members** (*group\_id, group\_urlname, max\_results, test=False*)

Get all of the Meetup members for a specified group.

**Parameters**

- **group\_id** (int) – The Meetup ID of the group.
- **group\_urlname** (str) – The URL name of the group.
- **max\_results** (int) – The maximum number of results to return per API query.
- **test** (bool) – For testing.

**Returns** A matchable list of Meetup members

**Return type** (list of dict)

**Members → Groups**

Start with Meetup members and end up with Meetup groups.

**exception NoMemberFound** (*member\_id*)

Bases: Exception

Exception should no member be found by the Meetup API

**get\_member\_details** (*member\_id, max\_results*)

Hit the Meetup API for details of a specified member

**Parameters**

- **member\_id** (str) – A Meetup member ID
- **max\_results** (int) – The maximum number of results with each API hit

**Returns** Meetup API response json.

**Return type** list of dict

**get\_member\_groups** (*member\_info*)

Extract the groups data from Meetup membership information.

**Parameters**

- **member\_id** (str) – A Meetup member ID

- **member\_info** (list of dict) – Meetup member API response json.

**Returns** List of unique member-group combinations

**Return type** list of dict

## Groups → Group details

Start with Meetup groups and end up with Meetup group details.

**exception NoGroupFound (group\_urlname)**

Bases: Exception

Exception should no group be found by the Meetup API

**get\_group\_details (group\_urlname, max\_results, avoid\_exception=True)**

Hit the Meetup API for the details of a specified groups. :param group\_urlname: A Meetup group urlname :type group\_urlname: str :param max\_results: Total number of results to return per API request. :type max\_results: int

**Returns** Meetup API response data

**Return type** (list of dict)

## Utils

Common tools between the different data collection points.

**get\_api\_key ()**

Get a random API key from those listed in the environmental variable MEETUP\_API\_KEYS.

**save\_sample (json\_data, filename, k)**

Dump a sample of k items from row-oriented JSON data json\_data into file with name filename.

**flatten\_data (list\_json\_data, keys, \*\*kwargs)**

Flatten nested JSON data from a list of JSON objects, by a list of desired keys. Each element in the keys may also be an ordered iterable of keys, such that subsequent keys describe a path through the JSON to desired value. For example in order to extract `key1` and `key3` from:

```
{'key': <some_value>, 'key2' : {'key3': <some_value>}}
```

one would specify keys as:

```
['key1', ('key2', 'key3')]
```

### Parameters

- **list\_json\_data (json)** – Row-orientated JSON data.
- **keys (list)** – Mixed list of either: individual str keys for data values
- **are not nested; or sublists of str, as described above.**  
(which) –
- **\*\*kwargs** – Any constants to include in every flattened row of the output.

**Returns** Flattened row-orientated JSON data.

**Return type** json

**get\_members\_by\_percentile**(*engine*, *perc*=10)

Get the number of meetup group members for a given percentile from the database.

**Parameters**

- **engine** – A SQL alchemy connectable.
- **perc** (*int*) – A percentile to evaluate.

**Returns** The number of members corresponding to this percentile.

**Return type** members (float)

**get\_core\_topics**(*engine*, *core\_categories*, *members\_limit*, *perc*=99)

Get the most frequent topics from a selection of meetup categories, from the database.

**Parameters**

- **engine** – A SQL alchemy connectable.
- **core\_categories** (*list*) – A list of category\_shortnames.
- **members\_limit** (*int*) – Minimum number of members required in a group for it to be considered.
- **perc** (*int*) – A percentile to evaluate the most frequent topics.

**Returns** The set of most frequent topics.

**Return type** topics (set)

### 1.1.2 Health data

Initially for our project with the Robert Woods Johnson Foundation (RWJF), these procedures outline the data collection of health-specific data.

#### Collect NIH

Extract all of the NIH World RePORTER data via their static data dump. N\_TABS outputs are produced in CSV format (concatenated across all years), where N\_TABS corresponds to the number of tabs in the main table found at:

[https://exporter.nih.gov/ExPORTER\\_Catalog.aspx](https://exporter.nih.gov/ExPORTER_Catalog.aspx)

The data is transferred to the Nesta intermediate data bucket.

**get\_data\_urls**(*tab\_index*)

Get all CSV URLs from the *tab\_index*'th tab of the main table found at :code:`TOP\_URL.

**Parameters** **tab\_index** (*int*) – Tab number (0-indexed) of table to extract CSV URLs from.

**Returns** Title of the tab in the table. hrefs (list): List of URLs pointing to data CSVs.

**Return type** title (str)

**iterrows**(*url*)

Yield rows from the CSV (found at URL *url*) as JSON (well, dict objects).

**Parameters** **url** (*str*) – The URL at which a zipped-up CSV is found.

**Yields** dict object, representing one row of the CSV.

## Process NIH

Data cleaning and processing procedures for the NIH World Reporter data. Specifically, a lat/lon is generated for each city/country; and the formatting of date fields is unified.

### 1.1.3 NLP Utils

Standard tools for aiding natural language processing.

#### Preprocess

Tools for preprocessing text.

**tokenize\_document** (*text*, *remove\_stops=False*)

Preprocess a whole raw document. :param text: Raw string of text. :type text: str :param remove\_stops: Flag to remove english stopwords :type remove\_stops: bool

**Returns** List of preprocessed and tokenized documents

**clean\_and\_tokenize** (*text*, *remove\_stops*)

Preprocess a raw string/sentence of text. :param text: Raw string of text. :type text: str :param remove\_stops: Flag to remove english stopwords :type remove\_stops: bool

**Returns** Preprocessed tokens.

**Return type** tokens (list, str)

**filter\_by\_idf** (*documents*, *lower\_idf\_limit*, *upper\_idf\_limit*)

Remove (from documents) terms which are in a range of IDF values.

#### Parameters

- **documents** (*list*) – Either a list of str or a list of list of str to be filtered.
- **lower\_idf\_limit** (*float*) – Lower percentile (between 0 and 100) on which to exclude terms by their IDF.
- **upper\_idf\_limit** (*float*) – Upper percentile (between 0 and 100) on which to exclude terms by their IDF.

**Returns** Filtered documents

### 1.1.4 Geo Utils

Tools for processing of geographical data, such as geocoding.

#### geocode

Tools for geocoding.

**geocode** (\*\**request\_kwargs*)

Geocoder using the Open Street Map Nominatim API.

If there are multiple results the first one is returned (they are ranked by importance). The API usage policy allows maximum 1 request per second and no multithreading: <https://operations.osmfoundation.org/policies/nominatim/>

**Parameters** **request\_kwargs** (*dict*) – Parameters for OSM API.

**Returns** JSON from API response.

**retry\_if\_not\_value\_error**(*exception*)

Forces retry to exit if a ValueError is returned. Supplied to the ‘retry\_on\_exception’ argument in the retry decorator.

**Parameters** **exception** (*Exception*) – the raised exception, to check

**Returns** False if a ValueError, else True

**Return type** (bool)

**geocode\_dataframe**(*df*)

A wrapper for the geocode function to process a supplied dataframe using the city and country.

**Parameters** **df** (*dataframe*) – a dataframe containing city and country fields.

**Returns** a dataframe with a ‘coordinates’ column appended.

**geocode\_batch\_dataframe**(*df*, *city*=‘city’, *country*=‘country’, *latitude*=‘latitude’, *longitude*=‘longitude’, *query\_method*=‘both’)

Geocodes a dataframe, first by supplying the city and country to the api, if this fails a second attempt is made supplying the combination using the q= method. The supplied dataframe df is returned with additional columns appended, containing the latitude and longitude as floats.

**Parameters**

- **df** (`pandas.DataFrame`) – input dataframe
- **city** (`str`) – name of the input column containing the city
- **country** (`str`) – name of the input column containing the country
- **latitude** (`str`) – name of the output column containing the latitude
- **longitude** (`str`) – name of the output column containing the longitude
- **query\_method** (`int`) – query methods to attempt: ‘city\_country\_only’: city and country only ‘query\_only’: q method only ‘both’: city, country with fallback to q method

**Returns** original dataframe with lat and lon appended as floats

**Return type** (`pandas.DataFrame`)

**generate\_composite\_key**(*city=None*, *country=None*)

Generates a composite key to use as the primary key for the geographic data.

**Parameters**

- **city** (`str`) – name of the city
- **country** (`str`) – name of the country

**Returns** composite key

**Return type** (str)

## 1.1.5 Format Utils

Tools for formatting data, such as dates.

## datetools

Tools for processing dates in data.

### `extract_year(date)`

Use search for 4 digits in a row to identify the year and return as YYYY-01-01.

**Parameters** `date (str)` – The full date string.

**Returns** integer

### `extract_date(date, date_format='%Y-%m-%d', return_date_object=False)`

Determine the date format, convert and return in YYYY-MM-DD format.

**Parameters** `date (str)` – the full date string.

**Returns** Formatted date string.

## 1.1.6 Decorators

### ratelimit

Apply rate limiting at a threshold per second

#### `ratelimit(max_per_second)`

**Parameters** `max_per_second (float)` – Number of permitted hits per second

### schema\_transform

Apply a field name transformation to a data output from the wrapped function, such that specified field names are transformed and unspecified fields are dropped. A valid file would be formatted as shown:

```
[{"tier_0": "bad_col", "tier_1": "good_col"}, {"tier_0": "another_bad_col", "tier_1": "another_good_col"}, ...]
```

where `tier_0` and `tier_1` correspond to `from_key` and `to_key` in the below documentation.

#### `load_transformer(filename, from_key, to_key)`

#### `schema_transform(filename, from_key, to_key)`

##### Parameters

- `filename (str)` – A record-oriented JSON file path mapping field names denoted by `from_key` and `to_key`.
- `from_key (str)` – The key in file indicated by `filename` which indicates the field name to transform.
- `to_key (str)` – The key in file indicated by `filename` which what the field name indicated by `from_key` will be transformed to.

**Returns** Data in the format it was originally passed to the wrapper in, with specified field names transformed and unspecified fields dropped.

#### `schema_transformer(data, *, filename, from_key, to_key, ignore=[])`

Function version of the schema\_transformer wrapper. :param data: the data requiring the schema transformation :type data: dataframe OR list of dicts :param filename: the path to the schema json file :type filename: str :param from\_key: tier level of the data :type from\_key: str :param to\_key: tier level to be applied to the data :type to\_key: str :param ignore: optional list of fields, eg ids or keys which shouldn't be dropped :type ignore: list

**Returns** supplied data with schema applied

## 1.2 Code auditing

Packages are only accepted if they satisfy our internal auditing procedure:

- **Common sense requirements:**
  - **Either:**
    - \* The code produces at least one data or model output; **or**
    - \* The code provides a service which abstracts away significant complexity.
  - There is one unit test for each function or method, which lasts no longer than about 1 minute.
  - Each data or model output is produced from a single function or method, as described in the `__main__` of a specified file.
  - Can the nearest programmer (or equivalent) checkout and execute your tests from scratch?
  - Will the code be used to perform non-project specific tasks?
  - Does the process perform a logical task or fulfil a logical purpose?
- **If the code requires productionising, it satisfies one of the following conditions:**
  - a) There is a non-trivial pipeline, which would benefit from formal productionising.
  - b) A procedure is foreseen to be reperformed for new contexts with atomic differences in run conditions.
  - c) The output is a service which requires a pipeline.
  - d) The process is a regular / longitudinal data collection.
- **Basic PEP8 and style requirements:**
  - Docstrings for every exposable class, method and function.
  - Usage in a README.rst or in Docstring at the top of the file.
  - CamelCase class names.
  - Underscore separation of all other variable, function and module names.
  - No glaring programming no-nos.
  - Never use `print`: opt for `logging` instead.
- **Bureaucratic requirements:**
  - A requirements file\*.
  - The README file specifies the operating system and python version.



# CHAPTER 2

---

## Production

---

Nesta's production system is based on [Luigi](#) pipelines, and are designed to be entirely run on AWS via the batch service. The main Luigi server runs on a persistent EC2 instance. Beyond the well documented Luigi code, the main features of the nesta production system are:

- `luigihacks.autobatch`, which facilitates a managed `Luigi.Task` which is split, batched and combined in a single step. Currently only synchronous jobs are accepted. Asynchronous jobs (where downstream `Luigi.Task` jobs can be triggered) are a part of a longer term plan.
- `scripts.nesta_prepare_batch` which zips up the batchable with the specified environmental files and ships it to AWS S3.
- `scripts.nesta_docker_build` which builds a specified docker environment and ships it to AWS ECS.

### 2.1 How to put code into production at nesta

If you're completely new, check out our [training slides](#). In short, the steps you should go through when building production code are to:

1. Audit the package code, required to pass all auditing tests
2. Understand what environment is required
3. Write a Dockerfile and docker launch script for this under `scripts/docker_recipes`
4. Build the Docker environment (run: `docker_build <recipe_name>` from any directory)
5. Build and test the batchable(s)
6. Build and test a Luigi pipeline
7. [...] Need to have steps here which estimate run time cost parameters. Could use `tests.py` to estimate this. [...]
8. Run the full chain

## 2.2 Code and scripts

### 2.2.1 Routines

All of our pipelines, implemented as Luigi routines. Some of these pipelines (at least partly) rely on batch computing (via AWS Batch), where the ‘batched’ scripts (*run.py* modules) are described in `core.batchables`. Other than `luighacks.autobatch`, which is respectively documented, the routine procedure follows the core [Luigi](#) documentation.

### Examples

Examples of Luigi routines, from which all other `nesta` production routines can be built. Currently we have examples of routines with S3 and database (MySQL) IO, and routines which are entirely batched.

We’d recommend reading Spotify’s [Luigi](#) documentation, and also checking the [Luigi Hacks](#) documentation which contains modified Luigi modules which (who knows) one day we will suggest as pull requests.

#### S3 Example

An example of building a pipeline with S3 Targets

```
class InputData(*args, **kwargs)
    Bases: luigi.task.ExternalTask
    Dummy task acting as the single input data source

    output()
        Points to the S3 Target

class SomeTask(*args, **kwargs)
    Bases: luigi.task.Task
    An intermediate task which increments the age of the muppets by 1 year.

        Parameters date (datetime) – Date used to label the outputs
        date = <luigi.parameter.DateParameter object>

        requires()
            Gets the input data (json containing muppet name and age)

        output()
            Points to the S3 Target

        run()
            Increments the muppets' ages by 1

class FinalTask(*args, **kwargs)
    Bases: luigi.task.Task
    The root task, which adds the surname ‘Muppet’ to the names of the muppets.

        Parameters date (datetime) – Date used to label the outputs
        date = <luigi.parameter.DateParameter object>

        requires()
            Get data from the intermediate task
```

---

```
output()
    Points to the S3 Target

run()
    Appends 'Muppet' the muppets' names
```

## Database example

An example of building a pipeline with database Targets

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

    Dummy task acting as the single input data source.

Parameters
    • date (datetime) – Date used to label the outputs
    • db_config – (dict) The input database configuration

date = <luigi.parameter.DateParameter object>
db_config = <luigi.parameter.DictParameter object>

output()
    Points to the input database target

run()
    Example of marking the update table
```

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

Task which increments the age of the muppets, by first selecting muppets with an age less than *max\_age*.

```
Parameters
    • date (datetime) – Date used to label the outputs
    • max_age (int) – Maximum age of muppets to select from the database
    • in_db_config – (dict) The input database configuration
    • out_db_config – (dict) The output database configuration

date = <luigi.parameter.DateParameter object>
max_age = <luigi.parameter.IntParameter object>
in_db_config = <luigi.parameter.DictParameter object>
out_db_config = <luigi.parameter.DictParameter object>

requires()
    Gets the input database engine

output()
    Points to the output database engine

run()
    Increments the muppets' ages by 1
```

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

A dummy root task, which collects the database configurations and executes the central task.

    Parameters date (datetime) – Date used to label the outputs
    date = <luigi.parameter.DateParameter object>
    requires()
        Collects the database configurations and executes the central task.
```

## Batch Example

An example of building a pipeline with batched tasks.

```
class SomeInitialTask(*args, **kwargs)
Bases: luigi.task.ExternalTask

Dummy task acting as the single input data source

output()
    Points to the S3 Target

class SomeBatchTask(*args, **kwargs)
Bases: nesta.core.luigihacks.autobatch.AutoBatchTask

A set of batched tasks which increments the age of the muppets by 1 year.
```

### Parameters

- **date** (datetime) – Date used to label the outputs
- **batchable** (str) – Path to the directory containing the run.py batchable
- **job\_def** (str) – Name of the AWS job definition
- **job\_name** (str) – Name given to this AWS batch job
- **job\_queue** (str) – AWS batch queue
- **region\_name** (str) – AWS region from which to batch
- **poll\_time** (int) – Time between querying the AWS batch job status

```
    date = <luigi.parameter.DateParameter object>
```

```
    requires()
```

Gets the input data (json containing muppet name and age)

```
    output()
```

Points to the S3 Target

```
    prepare()
```

Prepare the batch job parameters

```
    combine(job_params)
```

Combine the outputs from the batch jobs

```
class RootTask(*args, **kwargs)
```

Bases: luigi.task.Task

The root task, which adds the surname ‘Muppet’ to the names of the muppets.

Parameters **date** (datetime) – Date used to label the outputs

```

date = <luigi.parameter.DateParameter object>
requires()
    Get the output from the batchtask
output()
    Points to the S3 Target
run()
    Appends ‘Muppet’ the muppets’ names

```

## arXiv data (technical research)

Data collection and processing pipeline for arXiv data, principally for the [arXlive](#) platform. This pipeline orchestrates the collection of arXiv data, enrichment (via MAG and GRID), topic modelling, and novelty (lolvelty) measurement.

### Root task (arXlive)

Luigi routine to collect all data from the arXiv api and load it to MySQL, pipe to Elasticsearch, perform topic modelling, generate plots and measure novelty.

```

class RootTask (*args, **kwargs)
    Bases: luigi.task.WrapperTask
    A dummy root task, which collects the database configurations and executes the central task.

```

#### Parameters

- **date** (`datetime`) – Date used to label the outputs
- **db\_config\_path** (`str`) – Path to the MySQL database configuration
- **production** (`bool`) – Flag indicating whether running in testing mode (False, default), or production mode (True).
- **drop\_and\_recreate** (`bool`) – If in test mode, allows dropping the dev index from the ES database.

```

date = <luigi.parameter.DateParameter object>
db_config_path = <luigi.parameter.Parameter object>
production = <luigi.parameter.BoolParameter object>
drop_and_recreate = <luigi.parameter.BoolParameter object>
articles_from_date = <luigi.parameter.Parameter object>
insert_batch_size = <luigi.parameter.IntParameter object>
debug = <luigi.parameter.BoolParameter object>
requires()
    Collects the database configurations and executes the central task.

```

```

class EsOnlyRootTask (*args, **kwargs)
    Bases: luigi.task.WrapperTask
    A dummy root task, which collects the database configurations and executes the central task.

```

#### Parameters

- **date** (`datetime`) – Date used to label the outputs

- **db\_config\_path** (*str*) – Path to the MySQL database configuration
- **production** (*bool*) – Flag indicating whether running in testing mode (False, default), or production mode (True).
- **drop\_and\_recreate** (*bool*) – If in test mode, allows dropping the dev index from the ES database.

```
date = <luigi.parameter.DateParameter object>
db_config_path = <luigi.parameter.Parameter object>
production = <luigi.parameter.BoolParameter object>
drop_and_recreate = <luigi.parameter.BoolParameter object>
requires()
```

Collects the database configurations and executes the central task.

## Collection task

Luigi routine to collect new data from the arXiv api and load it to MySQL.

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

Collect new data from the arXiv api and dump the data in the MySQL server.

### Parameters

- **date** (*datetime*) – Datetime used to label the outputs
- **\_routine\_id** (*str*) – String used to label the AWS task
- **db\_config\_env** (*str*) – environmental variable pointing to the db config file
- **db\_config\_path** (*str*) – The output database configuration
- **insert\_batch\_size** (*int*) – number of records to insert into the database at once
- **articles\_from\_date** (*str*) – new and updated articles from this date will be retrieved. Must be in YYYY-MM-DD format

```
date = <luigi.parameter.DateParameter object>
test = <luigi.parameter.BoolParameter object>
db_config_env = <luigi.parameter.Parameter object>
db_config_path = <luigi.parameter.Parameter object>
insert_batch_size = <luigi.parameter.IntParameter object>
articles_from_date = <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

## Date task

Luigi wrapper to identify the date since the last iterative data collection

```
class DateTask (*args, **kwargs)
```

Bases: luigi.task.WrapperTask

Collect new data from the arXiv api and dump the data in the MySQL server.

### Parameters

- **date** (*datetime*) – Datetime used to label the outputs
- **\_routine\_id** (*str*) – String used to label the AWS task
- **db\_config\_env** (*str*) – environmental variable pointing to the db config file
- **db\_config\_path** (*str*) – The output database configuration
- **insert\_batch\_size** (*int*) – number of records to insert into the database at once
- **articles\_from\_date** (*str*) – new and updated articles from this date will be retrieved. Must be in YYYY-MM-DD format

```
date = <luigi.parameter.DateParameter object>
```

```
test = <luigi.parameter.BoolParameter object>
```

```
db_config_path = <luigi.parameter.Parameter object>
```

```
db_config_env = <luigi.parameter.Parameter object>
```

```
insert_batch_size = <luigi.parameter.IntParameter object>
```

```
articles_from_date = <luigi.parameter.Parameter object>
```

```
requires()
```

Collects the last date of successful update from the database and launches the iterative data collection task.

## arXiv enriched with MAG (API)

Luigi routine to query the Microsoft Academic Graph for additional data and append it to the exiting data in the database.

```
class QueryMagTask (*args, **kwargs)
```

Bases: luigi.task.Task

Query the MAG for additional data to append to the arxiv articles, primarily the fields of study.

### Parameters

- **date** (*datetime*) – Datetime used to label the outputs
- **\_routine\_id** (*str*) – String used to label the AWS task
- **db\_config\_env** (*str*) – environmental variable pointing to the db config file
- **db\_config\_path** (*str*) – The output database configuration
- **mag\_config\_path** (*str*) – Microsoft Academic Graph Api key configuration path
- **insert\_batch\_size** (*int*) – number of records to insert into the database at once (not used in this task but passed down to others)

- **articles\_from\_date** (*str*) – new and updated articles from this date will be retrieved. Must be in YYYY-MM-DD format (not used in this task but passed down to others)

```
date = <luigi.parameter.DateParameter object>
test = <luigi.parameter.BoolParameter object>
db_config_env = <luigi.parameter.Parameter object>
db_config_path = <luigi.parameter.Parameter object>
mag_config_path = <luigi.parameter.Parameter object>
insert_batch_size = <luigi.parameter.IntParameter object>
articles_from_date = <luigi.parameter.Parameter object>
```

**output()**

Points to the output database engine

**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

**run()**

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

See Task.run

## arXiv enriched with MAG (SPARQL)

Luigi routine to query the Microsoft Academic Graph for additional data and append it to the exiting data in the database. This is to collect information which is difficult to retrieve via the MAG API.

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

**Query the MAG for additional data to append to the arxiv articles,** primarily the fields of study.

### Parameters

- **date** (*datetime*) – Datetime used to label the outputs
- **\_routine\_id** (*str*) – String used to label the AWS task
- **db\_config\_env** (*str*) – environmental variable pointing to the db config file
- **db\_config\_path** (*str*) – The output database configuration
- **mag\_config\_path** (*str*) – Microsoft Academic Graph Api key configuration path
- **insert\_batch\_size** (*int*) – number of records to insert into the database at once (not used in this task but passed down to others)
- **articles\_from\_date** (*str*) – new and updated articles from this date will be retrieved. Must be in YYYY-MM-DD format (not used in this task but passed down to others)

```
date = <luigi.parameter.DateParameter object>
```

```

test = <luigi.parameter.BoolParameter object>
db_config_env = <luigi.parameter.Parameter object>
db_config_path = <luigi.parameter.Parameter object>
mag_config_path = <luigi.parameter.Parameter object>
insert_batch_size = <luigi.parameter.IntParameter object>
articles_from_date = <luigi.parameter.Parameter object>
output()
    Points to the output database engine

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

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

    See Task.run

```

## arXiv enriched with GRID

Luigi routine to lookup arXiv author's institutes via the GRID data, in order to "geocode" arXiv articles. The matching of institute name to GRID data is done via smart(ish) fuzzy matching, which then gives a confidence score per match.

```

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

    Join arxiv articles with GRID data for institute addresses and geocoding.

    Parameters
        • date (datetime) – Datetime used to label the outputs
        • _routine_id (str) – String used to label the AWS task
        • db_config_env (str) – environmental variable pointing to the db config file
        • db_config_path (str) – The output database configuration
        • mag_config_path (str) – Microsoft Academic Graph Api key configuration path
        • insert_batch_size (int) – number of records to insert into the database at once (not used in this task but passed down to others)
        • articles_from_date (str) – new and updated articles from this date will be retrieved. Must be in YYYY-MM-DD format (not used in this task but passed down to others)

date = <luigi.parameter.DateParameter object>
test = <luigi.parameter.BoolParameter object>
db_config_env = <luigi.parameter.Parameter object>
db_config_path = <luigi.parameter.Parameter object>
mag_config_path = <luigi.parameter.Parameter object>

```

```
insert_batch_size = <luigi.parameter.IntParameter object>
articles_from_date = <luigi.parameter.Parameter object>
output()
    Points to the output database engine
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
run()
    The task run method, to be overridden in a subclass.
See Task.run
```

## [AutoML] Topic modelling (CorEx)

Automated topic modelling of arXiv articles via the CorEx algorithm. See `topic_process_task_chain.json` for the full processing chain, but in brief: Vectorization is performed, followed by n-gramming (a lookup via Wiktionary) and then topics via CorEx.

```
class PrepareArxivS3Data(*args, **kwargs)
Bases: luigi.task.Task
Task that pipes SQL text fields to a number of S3 JSON files. This is particularly useful for preparing autoML tasks.
s3_path_out = <luigi.parameter.Parameter object>
db_conf_env = <luigi.parameter.Parameter object>
chunksize = <luigi.parameter.IntParameter object>
test = <luigi.parameter.BoolParameter object>
grid_task_kwargs = <nesta.core.luigihacks.parameter.DictParameterPlus 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
output()
    The output that this Task produces.
The output of the Task determines if the Task needs to be run—the task is considered finished iff the outputs all exist. Subclasses should override this method to return a single Target or a list of Target instances.
Implementation note If running multiple workers, the output must be a resource that is accessible by all workers, such as a DFS or database. Otherwise, workers might compute the same output since they don't see the work done by other workers.
See Task.output
```

```
write_to_s3(data, ichunk)

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

    See Task.run

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

        s3_path_prefix = <luigi.parameter.Parameter object>
        raw_s3_path_prefix = <luigi.parameter.Parameter object>
        data_path = <luigi.parameter.Parameter object>
        date = <luigi.parameter.DateParameter object>
        db_config_path = <luigi.parameter.Parameter object>
        db_conf_env = <luigi.parameter.Parameter object>
        test = <luigi.parameter.BoolParameter object>
        insert_batch_size = <luigi.parameter.IntParameter object>
        cherry_picked = <luigi.parameter.Parameter object>
        grid_task_kwargs = <nesta.core.luighacks.parameter.DictParameterPlus object>

        output()
            Points to the output database engine

        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

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

            See Task.run
```

## Pipe to elasticsearch

Luigi routine to load the Arxiv data from MYSQL into Elasticsearch.

```
class ArxivESTask(*args, **kwargs)
    Bases: nesta.core.luighacks.sql2estask.Sql2EsTask

        date = <luigi.parameter.DateParameter object>
        drop_and_recreate = <luigi.parameter.BoolParameter object>
        grid_task_kwargs = <nesta.core.luighacks.parameter.DictParameterPlus 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

## Elasticsearch tokenize

Tokenize arXiv field, which allows the search to leverage high quality n-grams, as provided by the Ngrammer module.

```
class ArxivESTokenTask(*args, **kwargs)
    Bases: nesta.core.luigihacks.estask.ElasticsearchTask

    done_ids()
        All document ids which do not require processing. If you want to avoid writing that function see
        LazyElasticsearchTask.

        Returns A set of document ids, not to be processed.

        Return type done_ids (set)

class ArxivTokenRootTask(*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

## Estimate novelty (lolvelty)

Estimate the novelty of each article via the lolvelty algorithm. This is performed on a document-by-document basis and is regrettably very slow since it is computationally very expensive for the Elasticsearch server.

```
class ArxivElasticsearchTask(*args, **kwargs)
    Bases: nesta.core.luigihacks.estask.ElasticsearchTask

    date = <luigi.parameter.DateParameter object>
    drop_and_recreate = <luigi.parameter.BoolParameter object>
    grid_task_kwargs = <nesta.core.luigihacks.parameter.DictParameterPlus object>
    done_ids()
        All document ids which do not require processing. If you want to avoid writing that function see
        LazyElasticsearchTask.

        Returns A set of document ids, not to be processed.

        Return type done_ids (set)
```

**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 ArxivLoveltyRootTask(*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

## “Deep learning, Deep Change” analysis

Luigi routine to perform the analysis from the Deep learning, deep change paper, placing the results in an S3 bucket to be picked up by the arXlive front end.

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

Extract and analyse arXiv data to produce data and charts for the arXlive front end to consume.

### Proposed charts:

1. distribution of dl/non dl papers by country (horizontal bar)
2. distribution of dl/non dl papers by city (horizontal bar)
3. % ML papers by year (line)
4. share of ML activity in arxiv subjects, pre/post 2012 (horizontal point / slope)
5. rca, pre/post 2012 by country (horizontal point / slope)
6. rca over time, citation > mean & top 50 countries (horizontal violin) [NOT DONE]

### Proposed table data:

1. top countries by rca (moving window of last 12 months?) [NOT DONE]

### Parameters

- **date** (*datetime*) – Datetime used to label the outputs
- **\_routine\_id** (*str*) – String used to label the AWS task
- **db\_config\_env** (*str*) – environmental variable pointing to the db config file
- **db\_config\_path** (*str*) – The output database configuration

- `mag_config_path` (`str`) – Microsoft Academic Graph Api key configuration path
- `insert_batch_size` (`int`) – number of records to insert into the database at once (not used in this task but passed down to others)
- `articles_from_date` (`str`) – new and updated articles from this date will be retrieved. Must be in YYYY-MM-DD format (not used in this task but passed down to others)

```
date = <luigi.parameter.DateParameter object>
test = <luigi.parameter.BoolParameter object>
db_config_env = <luigi.parameter.Parameter object>
db_config_path = <luigi.parameter.Parameter object>
mag_config_path = <luigi.parameter.Parameter object>
insert_batch_size = <luigi.parameter.IntParameter object>
articles_from_date = <luigi.parameter.Parameter object>
s3_path_prefix = <luigi.parameter.Parameter object>
raw_data_path = <luigi.parameter.Parameter object>
grid_task_kwargs = <nesta.core.luighacks.parameter.DictParameterPlus object>
cherry_picked = <luigi.parameter.Parameter object>
output()
    Points to the output database engine
```

**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

**run()**

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

See Task.run

```
class StandaloneAnalysisTask(*args, **kwargs)
    Bases: nestas.core.routines.arxiv.deepchange_analysis_task.AnalysisTask

    date = <luigi.parameter.DateParameter object>
    production = <luigi.parameter.BoolParameter object>
    test = <luigi.parameter.BoolParameter object>
    db_config_env = <luigi.parameter.Parameter object>
    db_config_path = <luigi.parameter.Parameter object>
    mag_config_path = <luigi.parameter.Parameter object>
    insert_batch_size = <luigi.parameter.IntParameter object>
    articles_from_date = <luigi.parameter.Parameter object>
    s3_path_prefix = <luigi.parameter.Parameter object>
```

```
raw_data_path = <luigi.parameter.Parameter object>
grid_task_kwargs = <nesta.core.luighacks.parameter.DictParameterPlus object>
cherry_picked = <luigi.parameter.Parameter 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

## CORDIS (EU funded research)

Generic pipeline (i.e. not project specific) to collect all CORDIS data, discovering all entities by crawling an unofficial API.

## H2020 and FP7 Data Collection

Collection of H2020 and FP7 projects, organisations, publications and topics from the unofficial API.

```
class CordisCollectTask(*args, **kwargs)
    Bases: nesta.core.luighacks.autobatch.AutoBatchTask

    process_batch_size = <luigi.parameter.IntParameter object>
    intermediate_bucket = <luigi.parameter.Parameter object>
    db_config_path = <luigi.parameter.Parameter object>
    db_config_env = <luigi.parameter.Parameter object>
    routine_id = <luigi.parameter.Parameter object>

    output()
        Points to the output database engine
```

**prepare()**

You should implement a method which returns a list of dict, where each dict corresponds to inputs to the batchable. Each row of the output must at least contain the following keys:

- **done (bool)**: indicating whether the job has already been finished.
- **outinfo (str)**: Text indicating e.g. the location of the output, for use in the batch job and for *combine* method

**Returns** list of dict

**combine (job\_params)**

You should implement a method which collects the outputs specified by the **outinfo** key of `job_params`, which is the output from the `prepare` method. This method should finally write to the luigi.Target output.

**Parameters** `job_params` (list of dict) – The batchable job parameters, as returned from the `prepare` method.

```
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

## Crunchbase (private sector companies)

**NB:** The Crunchbase pipeline may not work until [this issue](#) has been resolved.

Data collection and processing pipeline of Crunchbase data, principally for the [healthMosaic](#) platform.

### Root task (HealthMosaic)

Luigi routine to collect all data from the Crunchbase data dump and load it to MySQL, pipe to Elasticsearch, label projects as being health-related, assign mesh terms and deduplicate.

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

A dummy root task, which collects the database configurations and executes the central task.

#### Parameters

- **date** (`datetime`) – Date used to label the outputs
- **db\_config\_path** (`str`) – Path to the MySQL database configuration
- **production** (`bool`) – Flag indicating whether running in testing mode (False, default), or production mode (True).

```
date = <luigi.parameter.DateParameter object>
drop_and_recreate = <luigi.parameter.BoolParameter object>
production = <luigi.parameter.BoolParameter object>
insert_batch_size = <luigi.parameter.IntParameter object>
requires()
```

Collects the database configurations and executes the central task.

## Get organisations

Luigi routine to collect organisations from Crunchbase data exports and load the data into MySQL.

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

Download tar file of Organization csvs and load them into the MySQL server.

#### Parameters

- **\_routine\_id** (`str`) – String used to label the AWS task

- **db\_config\_path** – (str) The output database configuration

```
date = <luigi.parameter.DateParameter object>
test = <luigi.parameter.BoolParameter object>
insert_batch_size = <luigi.parameter.IntParameter object>
db_config_env = <luigi.parameter.Parameter object>
```

**output()**  
Points to the output database engine

**run()**  
Collect and process organizations, categories and long descriptions.

## Non-organisation collection

Luigi routine to collect non-organisation Crunchbase data exports and load the data into MySQL.

Organizations, category\_groups, org\_parents and organization\_descriptions should have already been processed; this task picks up all other files to be imported.

```
class NonOrgCollectTask(*args, **kwargs)
    Bases: nesta.core.luigihacks.autobatch.AutoBatchTask
```

Download tar file of csvs and load them into the MySQL server.

### Parameters

- **date (datetime)** – Datetime used to label the outputs
- **\_routine\_id (str)** – String used to label the AWS task
- **db\_config\_path** – (str) The output database configuration

```
date = <luigi.parameter.DateParameter object>
db_config_path = <luigi.parameter.Parameter object>
insert_batch_size = <luigi.parameter.IntParameter 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

**output()**  
Points to the output database engine

**prepare()**  
Prepare the batch job parameters

**combine (job\_params)**  
Touch the checkpoint

## Geocoding

Luigi routines to geocode the Organization, FundingRound, Investor, Ipo and People tables.

```
class OrgGeocodeTask (*args, **kwargs)
    Bases: nesta.core.luigihacks.batchgeocode.GeocodeBatchTask
    date = <luigi.parameter.DateParameter object>
    insert_batch_size = <luigi.parameter.IntParameter object>
    output()
        Points to the output database engine
    requires()
        Collects the database configurations and executes the central task.
    combine(job_params)
        Touch the checkpoint

class FundingRoundGeocodeTask (*args, **kwargs)
    Bases: nesta.core.luigihacks.batchgeocode.GeocodeBatchTask
    date = <luigi.parameter.DateParameter object>
    insert_batch_size = <luigi.parameter.IntParameter object>
    output()
        Points to the output database engine
    requires()
        Collects the database configurations and executes the central task.
    combine(job_params)
        Touch the checkpoint

class InvestorGeocodeTask (*args, **kwargs)
    Bases: nesta.core.luigihacks.batchgeocode.GeocodeBatchTask
    date = <luigi.parameter.DateParameter object>
    insert_batch_size = <luigi.parameter.IntParameter object>
    output()
        Points to the output database engine
    requires()
        Collects the database configurations and executes the central task.
    combine(job_params)
        Touch the checkpoint

class IpoGeocodeTask (*args, **kwargs)
    Bases: nesta.core.luigihacks.batchgeocode.GeocodeBatchTask
    date = <luigi.parameter.DateParameter object>
    insert_batch_size = <luigi.parameter.IntParameter object>
    output()
        Points to the output database engine
    requires()
        Collects the database configurations and executes the central task.
```

```

combine(job_params)
    Touch the checkpoint

class PeopleGeocodeTask(*args, **kwargs)
    Bases: nesta.core.luigihacks.batchgeocode.GeocodeBatchTask
        date = <luigi.parameter.DateParameter object>
        insert_batch_size = <luigi.parameter.IntParameter object>
        output()
            Points to the output database engine
        requires()
            Collects the database configurations and executes the central task.
        combine(job_params)
            Touch the checkpoint

```

## Organisation health labeling

Luigi routine to determine if crunchbase orgs are involved in health and apply a label to the data in MYSQL.

```

class HealthLabelTask(*args, **kwargs)
    Bases: luigi.task.Task
        Apply health labels to the organisation data in MYSQL.

Parameters
    • date (datetime) – Datetime used to label the outputs
    • _routine_id (str) – String used to label the AWS task
    • test (bool) – True if in test mode
    • insert_batch_size (int) – Number of rows to insert into the db in a batch
    • db_config_env (str) – The output database envariable
    • bucket (str) – S3 bucket where the models are stored
    • vectoriser_key (str) – S3 key for the vectoriser model
    • classifier_key (str) – S3 key for the classifier model

        date = <luigi.parameter.DateParameter object>
        test = <luigi.parameter.BoolParameter object>
        insert_batch_size = <luigi.parameter.IntParameter object>
        db_config_env = <luigi.parameter.Parameter object>
        bucket = <luigi.parameter.Parameter object>
        vectoriser_key = <luigi.parameter.Parameter object>
        classifier_key = <luigi.parameter.Parameter 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

**output ()**

Points to the output database engine

**run ()**

Apply health labels using model.

## Merge in parent organisations

This task picks up the missed org\_parents table from the Crunchbase data dump and combines this with organizations.

```
class ParentIdCollectTask(*args, **kwargs)
```

Bases: luigi.task.Task

Download tar file of csvs and append parent\_ids to the organizations table.

### Parameters

- **date** (`datetime`) – Datetime used to label the outputs
- **\_routine\_id** (`str`) – String used to label the AWS task
- **db\_config\_env** (`str`) – The output database envariable
- **db\_config\_path** (`str`) – The output database configuration
- **insert\_batch\_size** (`int`) – number of rows to insert into the db in a batch

```
date = <luigi.parameter.DateParameter object>
```

```
test = <luigi.parameter.BoolParameter object>
```

```
db_config_env = <luigi.parameter.Parameter object>
```

```
db_config_path = <luigi.parameter.Parameter object>
```

```
insert_batch_size = <luigi.parameter.IntParameter 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

**output ()**

Points to the output database engine

**run ()**

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

See Task.run

## Apply mesh terms

Collects and combines Mesh terms from S3 and descriptions from MySQL.

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

Collects and combines Mesh terms from S3, and descriptions from MYSQL.

Parameters

- date (str) – Date used to label the outputs
- _routine_id (str) – String used to label the AWS task
- db_config_path (str) – Path to the MySQL database configuration

date = <luigi.parameter.DateParameter object>
test = <luigi.parameter.BoolParameter object>
db_config_env = <luigi.parameter.Parameter object>
db_config_path = <luigi.parameter.Parameter object>
insert_batch_size = <luigi.parameter.IntParameter object>
requires()
    Collects the configurations and executes the previous task.

output()
    Points to the output database engine

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

See Task.run
```

## Pipe data to Elasticsearch

Luigi routine to load the Crunchbase data from MYSQL into Elasticsearch.

Not all data is copied: organizations, categories and locations only. The data is flattened and it is all stored in the same index.

```
class CrunchbaseSql2EsTask(*args, **kwargs)
Bases: nesta.core.luigihacks.autobatch.AutoBatchTask

Download tar file of csvs and load them into the MySQL server.

Parameters

- date (datetime) – Datetime used to label the outputs
- _routine_id (str) – String used to label the AWS task
- db_config_env (str) – The output database envariable
- process_batch_size (int) – Number of rows to process in a batch
- insert_batch_size (int) – Number of rows to insert into the db in a batch
- intermediate_bucket (str) – S3 bucket where the list of ids for each batch are written

date = <luigi.parameter.DateParameter object>
db_config_env = <luigi.parameter.Parameter object>
process_batch_size = <luigi.parameter.IntParameter object>
```

```
insert_batch_size = <luigi.parameter.IntParameter object>
intermediate_bucket = <luigi.parameter.Parameter object>
drop_and_recreate = <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

**output()**

Points to the output database engine

**prepare()**

You should implement a method which returns a list of dict, where each dict corresponds to inputs to the batchable. Each row of the output must at least contain the following keys:

- **done** (bool): indicating whether the job has already been finished.
- **outinfo** (str): Text indicating e.g. the location of the output, for use in the batch job and for *combine* method

**Returns** list of dict

**combine(job\_params)**

Touch the checkpoint

## Novelty score (lolvelty)

Apply “lolvelty” score to Crunchbase data (in Elasticsearch). Note: this is a slow procedure that is applied on a document-by-document basis.

```
class CrunchbaseLolveltyRootTask(*args, **kwargs)
```

Bases: luigi.task.WrapperTask

Apply Lolvelty score to crunchbase data.

### Parameters

- **production** (bool) – Running in full production mode?
- **index** (str) – Elasticsearch index to append Lolvelty score to.
- **date** (datetime) – Date for timestamping this routine.

```
production = <luigi.parameter.BoolParameter object>
```

```
index = <luigi.parameter.Parameter 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

## EURITO (piping data to Elasticsearch)

Pipeline specific to EURITO for piping existing data to Elasticsearch. A recent “EU” cut of patstat data is transferred from the “main” patstat database, to Nesta’s central database.

### Preprocess PATSTAT data

Select the EU subset of patstat, by doc family id. This is will significantly speed up transfer to ES.

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

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

    output ()
        Points to the output database engine

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

See Task.run

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

    date = <luigi.parameter.DateParameter object>
    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

## Root Task (EURITO)

Pipe data from MySQL to Elasticsearch, for use with clio-lite.

```
kwarg_maker (dataset, routine_id)

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

    process_batch_size = <luigi.parameter.IntParameter object>
    production = <luigi.parameter.BoolParameter object>
    date = <luigi.parameter.DateParameter object>
    drop_and_recreate = <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

## Gateway to Research (UK publicly funded research)

Generic pipeline (i.e. not project specific) to collect all GtR data, discovering all entities by crawling the official API. The routine then geocodes and loads data to MYSQL.

### Root Task (generic)

Luigi routine to collect all GtR data, geocode and load to MYSQL.

**class RootTask(\*args, \*\*kwargs)**

Bases: luigi.task.WrapperTask

A dummy root task, which collects the database configurations and executes the central task.

**Parameters date (datetime)** – Date used to label the outputs

**date = <luigi.parameter.DateParameter object>**

**page\_size = <luigi.parameter.IntParameter object>**

**production = <luigi.parameter.BoolParameter object>**

**requires()**

Collects the database configurations and executes the central task.

### Data collection

Discover all GtR data via the API.

**class GtrTask(\*args, \*\*kwargs)**

Bases: nesta.core.luigihacks.autobatch.AutoBatchTask

Get all GtR data

**date = <luigi.parameter.DateParameter object>**

**page\_size = <luigi.parameter.IntParameter object>**

**output()**

Points to the input database target

**prepare()**

Prepare the batch job parameters

**combine(job\_params)**

Combine the outputs from the batch jobs

**class GtrOnlyRootTask(\*args, \*\*kwargs)**

Bases: luigi.task.WrapperTask

A dummy root task, which collects the database configurations and executes the central task.

---

**Parameters**

- date** (*datetime*) – Date used to label the outputs
- date** = <luigi.parameter.DateParameter object>
- page\_size** = <luigi.parameter.IntParameter object>
- production** = <luigi.parameter.BoolParameter object>
- requires()**

Collects the database configurations and executes the central task.

## Geocode

Apply geocoding to the collected GtR data. Add country name, iso codes and continent.

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

Perform geocoding on the collected GtR organisations data

### Parameters

- **\_routine\_id** (*str*) – String used to label the AWS task
- **db\_config\_path** – (*str*) The output database configuration

**date** = <luigi.parameter.DateParameter object>  
**test** = <luigi.parameter.BoolParameter object>  
**db\_config\_env** = <luigi.parameter.Parameter object>  
**page\_size** = <luigi.parameter.IntParameter object>

**requires()**

Collects the database configurations and executes the central task.

**output()**

Points to the output database engine

**run()**

Collect and process organizations, categories and long descriptions.

## NiH (health research)

Data collection and processing pipeline of NiH data, principally for the [healthMosaic](#) platform.

### Root Task (HealthMosaic)

Luigi routine to collect NIH World RePORTER data via the World ExPORTER data dump. The routine transfers the data into the MySQL database before processing and indexing the data to ElasticSearch.

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

A dummy root task, which collects the database configurations and executes the central task.

### Parameters

- **date** (*datetime*) – Date used to label the outputs
- **db\_config\_path** (*str*) – Path to the MySQL database configuration

- **production** (*bool*) – Flag indicating whether running in testing mode (False, default), or production mode (True).

```
date = <luigi.parameter.DateParameter object>
db_config_path = <luigi.parameter.Parameter object>
production = <luigi.parameter.BoolParameter object>
drop_and_recreate = <luigi.parameter.BoolParameter object>
requires()
    Collects the database configurations and executes the central task.
```

## Data collection

Luigi routine to collect NIH World RePORTER data via the World ExPORTER data dump.

```
exists(_class, **kwargs)
class CollectTask(*args, **kwargs)
Bases: nesta.core.luigihacks.autobatch.AutoBatchTask
Scrape CSVs from the World ExPORTER site and dump the data in the MySQL server.
```

### Parameters

- **date** (*datetime*) – Datetime used to label the outputs
- **\_routine\_id** (*str*) – String used to label the AWS task
- **db\_config\_path** – (*str*) The output database configuration

```
date = <luigi.parameter.DateParameter object>
db_config_path = <luigi.parameter.Parameter object>
output()
    Points to the output database engine
prepare()
    Prepare the batch job parameters
combine(job_params)
    Touch the checkpoint
```

## Pipe to Elasticsearch

Transfers the data from the MySQL database to to ElasticSearch.

```
class ProcessTask(*args, **kwargs)
Bases: nesta.core.luigihacks.autobatch.AutoBatchTask
A dummy root task, which collects the database configurations and executes the central task.
```

### Parameters

- **date** (*str*) – Date used to label the outputs
- **\_routine\_id** (*str*) – String used to label the AWS task
- **db\_config\_path** (*str*) – Path to the MySQL database configuration

```
date = <luigi.parameter.DateParameter object>
```

```

db_config_path = <luigi.parameter.Parameter object>
drop_and_recreate = <luigi.parameter.BoolParameter object>
requires()
    Collects the database configurations and executes the central task.

output()
    Points to the input database target

batch_limits(query, batch_size)
    Determines first and last ids for a batch.

Parameters

- query (object) – orm query object
- batch_size (int) – rows of data in a batch

Returns first (int), last (int) application_ids

prepare()
    You should implement a method which returns a list of dict, where each dict corresponds to inputs to the batchable. Each row of the output must at least contain the following keys:


- done (bool): indicating whether the job has already been finished.
- outinfo (str): Text indicating e.g. the location of the output, for use in the batch job and for combine method

Returns list of dict

combine(job_params)
    You should implement a method which collects the outputs specified by the outinfo key of job_params, which is the output from the prepare method. This method should finally write to the luigi.Target output.

Parameters job_params (list of dict) – The batchable job parameters, as returned from the prepare method.



---


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

    A dummy root task, which collects the database configurations and executes the central task.

Parameters

- date (datetime) – Date used to label the outputs
- db_config_path (str) – Path to the MySQL database configuration
- production (bool) – Flag indicating whether running in testing mode (False, default), or production mode (True).

date = <luigi.parameter.DateParameter object>
db_config_path = <luigi.parameter.Parameter object>
production = <luigi.parameter.BoolParameter object>
drop_and_recreate = <luigi.parameter.BoolParameter object>
requires()
    Collects the database configurations and executes the central task.

```

## Assign MeSH terms to abstracts

Assign MeSH terms to projects which have abstracts (note: not all projects have abstracts).

```
split_mesh_file_key(key)
subset_keys(es, es_config, keys)
class AbstractsMeshTask(*args, **kwargs)
Bases: nesta.core.luigihacks.autobatch.AutoBatchTask
```

Collects and combines Mesh terms from S3, Abstracts from MYSQL and projects in Elasticsearch.

### Parameters

- **date** (*str*) – Date used to label the outputs
- **\_routine\_id** (*str*) – String used to label the AWS task
- **db\_config\_path** (*str*) – Path to the MySQL database configuration

```
date = <luigi.parameter.DateParameter object>
db_config_path = <luigi.parameter.Parameter object>
drop_and_recreate = <luigi.parameter.BoolParameter object>
```

### requires()

Collects the configurations and executes the previous task.

### output()

Points to the input database target

### static get\_abstract\_file\_keys(bucket, key\_prefix)

Retrieves keys to meshed files from s3.

### Parameters

- **bucket** (*str*) – s3 bucket
- **key\_prefix** (*str*) – prefix to identify the files, ie the folder and start of a filename

### Returns

keys of the files

### Return type

(set of str)

### done\_check(es\_client, index, doc\_type, key)

Checks elasticsearch for mesh terms in the first and last documents in the batch.

### Parameters

- **es\_client** (*object*) – instantiated elasticsearch client
- **index** (*str*) – name of the index
- **doc\_type** (*str*) – name of the document type
- **key** (*str*) – filepath in s3

### Returns

True if both existing, otherwise False

### Return type

(bool)

### prepare()

You should implement a method which returns a list of dict, where each dict corresponds to inputs to the batchable. Each row of the output must at least contain the following keys:

- **done** (*bool*): indicating whether the job has already been finished.

- **outinfo (str)**: Text indicating e.g. the location of the output, for use in the batch job and for *combine* method

**Returns** list of dict

#### combine (job\_params)

You should implement a method which collects the outputs specified by the **outinfo** key of `job_params`, which is the output from the `prepare` method. This method should finally write to the `luigi.Target` output.

**Parameters** `job_params` (list of dict) – The batchable job parameters, as returned from the `prepare` method.

## Deduplication of near duplicates

Remove nears duplicates of projects from the data. Numeric fields (such as funding) are aggregated together.

### class DedupeTask (\*args, \*\*kwargs)

Bases: `nesta.core.luighacks.autobatch.AutoBatchTask`

`date` = <`luigi.parameter.DateParameter` object>

`routine_id` = <`luigi.parameter.Parameter` object>

`intermediate_bucket` = <`luigi.parameter.Parameter` object>

`db_config_path` = <`luigi.parameter.Parameter` object>

`process_batch_size` = <`luigi.parameter.IntParameter` object>

`drop_and_recreate` = <`luigi.parameter.BoolParameter` object>

`output ()`

Points to the output database engine

`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`

`prepare ()`

You should implement a method which returns a list of dict, where each dict corresponds to inputs to the batchable. Each row of the output must at least contain the following keys:

- **done (bool)**: indicating whether the job has already been finished.
- **outinfo (str)**: Text indicating e.g. the location of the output, for use in the batch job and for *combine* method

**Returns** list of dict

#### combine (job\_params)

Touch the checkpoint

### class NiHLoveltyRootTask (\*args, \*\*kwargs)

Bases: `luigi.task.WrapperTask`

```
production = <luigi.parameter.BoolParameter object>
index = <luigi.parameter.Parameter 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

## Meetup (social networking data)

**NB: The Meetup pipeline will not work until this issue has been resolved.**

Data collection and processing pipeline of Meetup data, principally for the [healthMosaic](#) platform.

### Collect data

Starting with a seed country (and Meetup category), extract all groups in that country and subsequently find all groups associated with all members of the original set of groups.

**chunks** (*l, n*)

Yield successive *n*-sized chunks from *l*.

**class CountryGroupsTask (\*args, \*\*kwargs)**

Bases: `nesta.core.luigihacks.autobatch.AutoBatchTask`

Extract all groups with corresponding category for this country.

Args:

`iso2 = <luigi.parameter.Parameter object>`

`category = <luigi.parameter.Parameter object>`

**output ()**

Points to the input database target

**prepare ()**

Prepare the batch job parameters

**combine (job\_params)**

Combine the outputs from the batch jobs

**class GroupsMembersTask (\*args, \*\*kwargs)**

Bases: `nesta.core.luigihacks.autobatch.AutoBatchTask`

#### Parameters

- **date** (`datetime`) – Date used to label the outputs
- **batchable** (`str`) – Path to the directory containing the `run.py` batchable
- **job\_def** (`str`) – Name of the AWS job definition
- **job\_name** (`str`) – Name given to this AWS batch job
- **job\_queue** (`str`) – AWS batch queue

- **region\_name** (*str*) – AWS region from which to batch
- **poll\_time** (*int*) – Time between querying the AWS batch job status

```
iso2 = <luigi.parameter.Parameter object>
category = <luigi.parameter.Parameter object>
requires()
    Gets the input data
output()
    Points to the DB target
prepare()
    Prepare the batch job parameters
combine(job_params)
    Combine the outputs from the batch jobs
```

**class MembersGroupsTask(\*args, \*\*kwargs)**  
Bases: `nesta.core.luigihacks.autobatch.AutoBatchTask`

**Parameters**

- **date** (*datetime*) – Date used to label the outputs
- **batchable** (*str*) – Path to the directory containing the run.py batchable
- **job\_def** (*str*) – Name of the AWS job definition
- **job\_name** (*str*) – Name given to this AWS batch job
- **job\_queue** (*str*) – AWS batch queue
- **region\_name** (*str*) – AWS region from which to batch
- **poll\_time** (*int*) – Time between querying the AWS batch job status

```
iso2 = <luigi.parameter.Parameter object>
category = <luigi.parameter.Parameter object>
requires()
    Gets the input data
output()
    Points to the DB target
prepare()
    Prepare the batch job parameters
combine(job_params)
    Combine the outputs from the batch jobs
```

**class GroupDetailsTask(\*args, \*\*kwargs)**  
Bases: `nesta.core.luigihacks.autobatch.AutoBatchTask`

The root task, which adds the surname ‘Muppet’ to the names of the muppets.

**Parameters** **date** (*datetime*) – Date used to label the outputs

```
iso2 = <luigi.parameter.Parameter object>
category = <luigi.parameter.Parameter object>
requires()
    Get the output from the batchtask
```

```
output()
    Points to the DB target

prepare()
    Prepare the batch job parameters

combine(job_params)
    Combine the outputs from the batch jobs

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

    A dummy root task, which collects the database configurations and executes the central task.

    Parameters date(datetime) – Date used to label the outputs
    date = <luigi.parameter.DateParameter object>
    iso2 = <luigi.parameter.Parameter object>
    category = <luigi.parameter.Parameter object>
    production = <luigi.parameter.BoolParameter object>

    requires()
        Collects the database configurations and executes the central task.
```

## Topic discovery

Task to automatically discover relevant topics from meetup data, defined as the most frequently occurring from a set of categories.

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

    Task to automatically discover relevant topics from meetup data, defined as the most frequently occurring from a set of categories.
```

### Parameters

- **db\_config\_env(str)** – Environmental variable pointing to the path of the DB config.
- **routine\_id(str)** – The routine UID.
- **core\_categories(list)** – A list of category\_shortnames from which to identify topics.
- **members\_perc(int)** – A percentile to evaluate the minimum number of members.
- **topic\_perc(int)** – A percentile to evaluate the most frequent topics.
- **test(bool)** – Test mode.

```
db_config_env = <luigi.parameter.Parameter object>
routine_id = <luigi.parameter.Parameter object>
core_categories = <luigi.parameter.ListParameter object>
members_perc = <luigi.parameter.IntParameter object>
topic_perc = <luigi.parameter.IntParameter object>
test = <luigi.parameter.BoolParameter object>
```

---

```
output()
    Points to the S3 Target

run()
    Extract the topics of interest
```

## Pipe to elasticsearch

Luigi routine to load the Meetup Group data from MySQL into Elasticsearch.

```
class MeetupHealthSql2EsTask(*args, **kwargs)
    Bases: nesta.core.luigihacks.sql2estask.Sql2EsTask

    Task to pipe meetup data to ES. For other arguments, see Sql2EsTask.
```

### Parameters

- **core\_categories (list)** – A list of category\_shortnames from which to identify topics.
- **members\_perc (int)** – A percentile to evaluate the minimum number of members.
- **topic\_perc (int)** – A percentile to evaluate the most frequent topics.

```
core_categories = <luigi.parameter.ListParameter object>
members_perc = <luigi.parameter.IntParameter object>
topic_perc = <luigi.parameter.IntParameter 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>
    core_categories = <luigi.parameter.ListParameter object>
    members_perc = <luigi.parameter.IntParameter object>
    topic_perc = <luigi.parameter.IntParameter object>
    db_config_env = <luigi.parameter.Parameter object>
    drop_and_recreate = <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

## Novelty score (lolvelty)

Apply “lolvelty” score to Meetup data.

```
class MeetupLoveltyRootTask(*args, **kwargs)
```

Bases: luigi.task.WrapperTask

Apply Lovelty score to meetup data.

### Parameters

- **production** (*bool*) – Running in full production mode?
- **index** (*str*) – Elasticsearch index to append Lovelty score to.
- **date** (*datetime*) – Date for timestamping this routine.

```
production = <luigi.parameter.BoolParameter object>
```

```
index = <luigi.parameter.Parameter 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

## 2.2.2 Batchables

Packets of code to be batched by `core.routines` routines. Each packet should sit in it’s own directory, with a file called `run.py`, containing a ‘main’ function called `run()` which will be executed on the AWS batch system.

Each `run.py` should expect an environment parameter called `BATCHPAR_outfile` which should provide information on the output location. Other input parameters should be prefixed with `BATCHPAR_`, as set in `core.routines` routine.

### Data / project specific batchables

#### Example

There are two batchable examples listed here. The first is a module which will be run if you try executing the `batch_example` luigi routine. The second is purely meant as a template, if you are learning the design pattern for `nesta`’s luigi batchables.

#### run.py (batch\_example)

The batchable for the `routines.examples.batch_example`, which simply increments a muppet’s age by one unit.

```
run()
```

Gets the name and age of the muppet, and increments the age. The result is transferred to S3.

## run.py (template\_batchable)

This is a pretty generic example of how your run.py might look. It reads and writes from a table, and hits the S3 “checkpoint” at the end.

```
run()
```

## arXiv data (technical research)

### run.py (arxiv\_elasticsearch)

Transfer arXiv data from MySQL to elasticsearch, principally intended for the arXlive hierarXy searchkit front-end.

**hierarchy\_field**(row\_data)

Generate hierarchical representation of fields of study, which is required for the searchkit interface.

```
run()
```

## CORDIS (EU-funded research)

### run.py (cordis\_api)

Transfer data on organisations, projects and outputs from the Cordis API on a project-by-project basis.

**extract\_core\_orgs**(orgs, project\_rcn)

Separate a project-organisation (which) is likely to be a department, with a non-unique address.

#### Parameters

- **orgs** (*list*) – List of organisations to process (NB: this will be modified)
- **project\_rcn** (*str*) – The record number of this project

**Returns** The unique ‘parent’ organisations.

**Return type** core\_orgs (*list*)

**prepare\_data**(items, rcn)

Append the project code (‘RCN’) to each “row” (dict) of data (list)

**split\_links**(items, project\_rcn)

Generate link table items for each item (dict) in items (list) for the project

```
run()
```

## Crunchbase data (private companies)

**NB: The Crunchbase pipeline may not work until this issue has been resolved.**

Batchables for the collection and processing of Crunchbase data. As documented under *packages* and *routines*, the pipeline is executed in the following order (documentation for the *run.py* files is given below, which isn’t super-informative. You’re better off looking under packages and routines).

The data is collected from proprietary data dumps, parsed into MySQL (tier 0) and then piped into Elasticsearch (tier 1), post-processing.

### **run.py (crunchbase\_collect)**

Collect Crunchbase data from the proprietary data dump and pipe into the MySQL database.

```
run()
```

### **run.py (crunchbase\_elasticsearch)**

Pipe Crunchbase data from MySQL to Elasticsearch.

```
run()
```

## **EURITO**

Batchables for processing data (which has already been collected elsewhere within this codebase) for the [EURITO](#) project. All of these batchables pipe the data into an Elasticsearch database, which is then cloned by EURITO.

### **run.py (arxiv\_eu)**

Transfer pre-collected arXiv data from MySQL to Elasticsearch, whilst labelling arXiv articles as being EU or not. This differs slightly from the [arXlive](#) pipeline, by reflecting the EURITO project more specifically, and allowing more in depth analysis of MAG fields of study.

```
run()
```

### **run.py (crunchbase\_eu)**

Transfer pre-collected Crunchbase data from MySQL to Elasticsearch.

```
run()
```

### **run.py (patstat\_eu)**

Transfer pre-collected PATSTAT data from MySQL to Elasticsearch. Only EU patents since the year 2000 are considered. The patents are grouped by patent families.

```
select_text(objs, lang_field, text_field)
metadata(orm, session, appln_ids, field_selector=None)
run()
```

## **GtR (UK publicly funded research)**

Batchable tools for collecting and processing GtR data. As documented under packages and routines, the pipeline is executed in the following order (documentation for the run.py files is given below, which isn't super-informative. You're better off looking under packages and routines).

The data is collected by traversing the graph exposed by the GtR API, and is parsed into MySQL (tier 0). There is a further module for directly generating document embeddings of GtR project descriptions, which can be used for finding topics.

## run.py (collect\_gtr)

Starting from GtR projects, iteratively and recursively discover all GtR entities by crawling the API.

`run()`

## run.py (embed\_topics)

Document embedding of GtR data. Would be better if this was generalized (i.e. not GtR specific), and migrated to batchables.nlp [see <https://github.com/nestauk/nesta/issues/203>]

`run()`

## NiH data (health research)

Batchables for the collection and processing of NiH data. As documented under *packages* and *routines*, the pipeline is executed in the following order (documentation for the *run.py* files is given below, which isn't super-informative. You're better off looking under packages and routines).

The data is collected from official data dumps, parsed into MySQL (tier 0) and then piped into Elasticsearch (tier 1), post-processing.

### run.py (nih\_collect\_data)

Collect NiH table from the official data dump, based on the name of the table. The data is piped into the MySQL database.

`run()`

### run.py (nih\_process\_data)

Geocode NiH data (from MySQL) and pipe into Elasticsearch.

`run()`

### run.py (nih\_abstract\_mesh\_data)

Retrieve NiH abstracts from MySQL, assign pre-calculated MeSH terms for each abstract, and pipe data into Elasticsearch. Exact abstract duplicates are removed at this stage.

#### `clean_abstract` (*abstract*)

Removes multiple spaces, tabs and newlines.

**Parameters** `abstract` (*str*) – text to be cleaned

**Returns** (*str*): cleaned text

`run()`

## run.py (nih\_dedupe)

Deduplicate NiH articles based on similarity scores using Elasticsearch's document similarity API. Similarity is calculated based on the description of the project, the project abstract and the title of the project. Funding information is aggregated (summed) across all deduplicated articles, for the total and annuals funds.

`get_value (obj, key)`

Retrieve a value by key if exists, else return None.

`extract_yearly_funds (src)`

Extract yearly funds

`run ()`

## Meetup (social networking / knowledge exchange)

**NB: The meetup pipeline will not work until this issue has been resolved.**

Batchables for the Meetup data collection pipeline. As documented under *packages* and *routines*, the pipeline is executed in the following order (documentation for the *run.py* files is given below, which isn't super-informative. You're better off looking under packages and routines).

The `topic_tag_elasticsearch` module is responsible for piping data to elasticsearch, as well as apply topic tags and filtering small groups out of the data.

## run.py (country\_groups)

Batchable for expanding from countries to groups

`run ()`

## run.py (groups\_members)

Batchable for expanding group members

`run ()`

## run.py (members\_groups)

Batchable for expanding from members to groups.

`run ()`

## run.py (group\_details)

Batchable for expanding group details

`run ()`

## **run.py (topic\_tag\_elasticsearch)**

Batchable for piping data to Elasticsearch, whilst implementing topic tags, and filtering groups with too few members (given by the 10th percentile of group size, to avoid “junk” groups).

`run()`

### **General-purpose batchables**

#### **Bulk geocoding**

##### **run.py (batch\_geocode)**

Geocode any row-delimited json data, with columns corresponding to a city/town/etc and country.

`run()`

### **Natural Language Processing**

Batchable utilities for NLP. Note that modules prefixed with [AUTOML] are designed to be launched by AutoMLTask, and those with the addition \* (i.e. [AUTOML\*]) are the designed to be the final task in an AutoMLTask chain (i.e. they provide a ‘loss’).

#### **[AutoML\*] run.py (corex\_topic\_model)**

Generate topics based on the CorEx algorithm. Loss is calculated from the total correlation explained.

`run()`

#### **[AutoML] run.py (ngrammer)**

Find and replace ngrams in a body of text, based on Wiktionary N-Grams. Whilst at it, the ngrammer also tokenizes and removes stop words (unless they occur within an n-gram)

`run()`

#### **[AutoML] run.py (tfidf)**

Applies TFIDF cuts to a dataset via environmental variables lower\_tfidf\_percentile and upper\_tfidf\_percentile.

##### **chunker (\_transformed, n\_chunks)**

Yield chunks from a numpy array.

##### **Parameters**

- **\_transformed** (`np.array`) – Array to split into chunks.
- **n\_chunks** (`int`) – Number of chunks to split the array into.

**Yields** chunk (`np.array`)

`run()`

## [AutoML] vectorizer (run.py)

Vectorizes (counts or binary) text data, and applies basic filtering of extreme term/document frequencies.

**term\_counts** (*dct, row, binary=False*)

Convert a single single document to term counts via a gensim dictionary.

### Parameters

- **dct** (*Dictionary*) – Gensim dictionary.
- **row** (*str*) – A document.
- **binary** (*bool*) – Binary rather than total count?

**Returns** dict of term id (from the Dictionary) to term count.

**optional** (*name, default*)

Defines optional env fields with default values

**merge\_lists** (*list\_of\_lists*)

Join a lists of lists into a single list. Returns an empty list if the input is not a list, which is expected to happen (from the ngrammer) if no long text was found

**run()**

## Novelty

Batchables for calculating measures of “novelty”.

### 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()**

## 2.2.3 ORMs

SQLAlchemy ORMs for the routines, which allows easy integration of testing (including automatic setup of test databases and tables).

### Meetup

```
class Group(**kwargs)
    Bases: sqlalchemy.ext.declarative.api.Base
    id
    urlname
    category_name
    category_shortname
    city
    country
```

---

```

created
description
lat
lon
members
name
topics
category_id
country_name
timestamp

class GroupMember(**kwargs)
Bases: sqlalchemy.ext.declarative.api.Base

Note: no foreign key constraint, since unknown groups will be found in the member expansion phase

group_id
group_urlname
member_id

```

## NIH schema

The schema for the World RePORTER data. Note that the schema was automatically generated based on an IPython hack session using the limits and properties of the data from the World ExPORTER (which explains the unusual VARCHAR limits).

```

class Projects(**kwargs)
Bases: sqlalchemy.ext.declarative.api.Base

application_id
activity
administering_ic
application_type
arra_funded
award_notice_date
budget_start
budget_end
cfda_code
core_project_num
ed_inst_type
foa_number
full_project_num
funding_ics

```

```
funding_mechanism
fy
ic_name
org_city
org_country
org_dept
org_district
org_duns
org_fips
org_ipf_code
org_name
org_state
org_zipcode
phr
pi_ids
pi_names
program_officer_name
project_start
project_end
project_terms
project_title
serial_number
study_section
study_section_name
suffix
support_year
direct_cost_amt
indirect_cost_amt
total_cost
subproject_id
total_cost_sub_project
nih_spending_cats

class Abstracts(**kwargs)
Bases: sqlalchemy.ext.declarative.api.Base
application_id
abstract_text
```

```

class Publications(**kwargs)
    Bases: sqlalchemy.ext.declarative.api.Base

        pmid
        author_name
        affiliation
        author_list
        country
        issn
        journal_issue
        journal_title
        journal_title_abbr
        journal_volume
        lang
        page_number
        pub_date
        pub_title
        pub_year
        pmc_id

class Patents(**kwargs)
    Bases: sqlalchemy.ext.declarative.api.Base

        patent_id
        patent_title
        project_id
        patent_org_name

```

## 2.2.4 Ontologies and schemas

### 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.

### 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*

## Tier 2

Although not-yet-implemented, the tier 2 schema is reserved for future graph ontologies. Don't expect any changes any time soon!

### 2.2.5 Luigi Hacks

Modifications and possible future contributions to the Luigi module.

#### autobatch

Automatic preparation, submission and consolidation of AWS batch tasks; as a single Luigi Task.

**command\_line** (*command*, *verbose=False*)

Execute command line tasks and return the final output line. This is particularly useful for extracting the AWS access keys directly from the OS; as well as executing the environment preparation script (core/scripts/nesta\_prepare\_batch.sh).

**class AutoBatchTask** (\*args, \*\*kwargs)  
Bases: luigi.task.Task, abc.ABC

A base class for automatically preparing and submitting AWS batch tasks.

Unlike regular Luigi Tasks, which require the user to override the `requires`, `output` and `run` methods, `AutoBatchTask` instead effectively replaces `run` with two new abstract methods: `prepare` and `combine`, which are repectively documented. With these abstract methods specified, `AutoBatchTask` will automatically prepare, submit, and combine one batch task (specified in `core.batchables`) per parameter set specified in the `prepare` method. The `combine` method will subsequently combine the outputs from the batch task.

#### Parameters

- **batchable** (*str*) – Path to the directory containing the run.py batchable
- **job\_def** (*str*) – Name of the AWS job definition
- **job\_name** (*str*) – Name given to this AWS batch job
- **job\_queue** (*str*) – AWS batch queue
- **region\_name** (*str*) – AWS region from which to batch
- **env\_files** (*list of str, optional*) – List of names pointing to local environmental files (for example local imports or scripts) which should be zipped up with the AWS batch job environment. Defaults to [].
- **vcpus** (*int, optional*) – Number of CPUs to request for the AWS batch job. Defaults to 1.

- **memory** (*int, optional*) – Memory to request for the AWS batch job. Defaults to 512 MiB.
- **max\_runs** (*int, optional*) – Number of batch jobs to run, which is useful for testing a subset of the full pipeline, or making cost predictions for AWS computing time. Defaults to *None*, implying that all jobs should be run.
- **poll\_time** (*int, optional*) – Time in seconds between querying the AWS batch job status. Defaults to 60.
- **success\_rate** (*float, optional*) – If the fraction of FAILED jobs exceeds success\_rate then the entire Task, along with any submitted AWS batch jobs, is killed. The fraction is calculated with respect to any jobs with RUNNING, SUCCEEDED or FAILED status. Defaults to 0.75.

```

batchable = <luigi.parameter.Parameter object>
job_def = <luigi.parameter.Parameter object>
job_name = <luigi.parameter.Parameter object>
job_queue = <luigi.parameter.Parameter object>
region_name = <luigi.parameter.Parameter object>
env_files = <luigi.parameter.ListParameter object>
vcpus = <luigi.parameter.IntParameter object>
memory = <luigi.parameter.IntParameter object>
max_runs = <luigi.parameter.IntParameter object>
timeout = <luigi.parameter.IntParameter object>
poll_time = <luigi.parameter.IntParameter object>
success_rate = <luigi.parameter.FloatParameter object>
test = <luigi.parameter.BoolParameter object>
max_live_jobs = <luigi.parameter.IntParameter object>
worker_timeout = inf
run()

```

DO NOT OVERRIDE THIS METHOD.

An implementation of the Luigi.Task.run method which is a wrapper around the prepare, execute and combine methods. Instead of overriding this method, you should implement prepare and combine methods in your class.

#### `prepare()`

You should implement a method which returns a list of dict, where each dict corresponds to inputs to the batchable. Each row of the output must at least contain the following keys:

- **done** (*bool*): indicating whether the job has already been finished.
- **outinfo** (*str*): Text indicating e.g. the location of the output, for use in the batch job and for *combine* method

**Returns** list of dict

**combine** (*job\_params*)

You should implement a method which collects the outputs specified by the **outinfo** key of *job\_params*, which is the output from the `prepare` method. This method should finally write to the `luigi.Target` output.

**Parameters** **job\_params** (list of dict) – The batchable job parameters, as returned from the `prepare` method.

**execute** (*job\_params*, *s3file\_timestamp*)

The secret sauce, which automatically submits and monitors the AWS batch jobs. Your AWS access key and id are automatically retrieved via the AWS CLI.

**Parameters**

- **job\_params** (list of dict) – The batchable job parameters, as returned from the `prepare` method. Each job is submitted from every item in this list. Each *dict* key-value pair is converted into an environmental variable in the batch job, with the variable name formed from the key, prefixed by `BATCHPAR_`.
- **s3file\_timestamp** (str) – The timestamp of the batchable zip file to be found on S3 by the AWS batch job.

## batchclient

NOTE: overwhelmingly based on [this](#), where the following documentation has been directly lifted. The main difference to the latter, is that AWS jobs are submitted via `**kwargs` in order to allow more flexibility (and probably more future-proofing if new parameters are added to `boto3`).

AWS Batch wrapper for Luigi

From the AWS website:

AWS Batch enables you to run batch computing workloads on the AWS Cloud.

Batch computing is a common way for developers, scientists, and engineers to access large amounts of compute resources, and AWS Batch removes the undifferentiated heavy lifting of configuring and managing the required infrastructure. AWS Batch is similar to traditional batch computing software. This service can efficiently provision resources in response to jobs submitted in order to eliminate capacity constraints, reduce compute costs, and deliver results quickly.

See [AWS Batch User Guide](#) for more details.

To use AWS Batch, you create a `jobDefinition` JSON that defines a `docker run` command, and then submit this JSON to the API to queue up the task. Behind the scenes, AWS Batch auto-scales a fleet of EC2 Container Service instances, monitors the load on these instances, and schedules the jobs.

This `boto3-powered` wrapper allows you to create Luigi Tasks to submit Batch `jobDefinition`'s`. You can either pass a dict (mapping directly to the ```jobDefinition` JSON) OR an Amazon Resource Name (arn) for a previously registered `jobDefinition`.

Requires:

- `boto3` package
- Amazon AWS credentials discoverable by `boto3` (e.g., by using `aws configure` from `awscli`)
- An enabled AWS Batch job queue configured to run on a compute environment.

Written and maintained by Jake Feala (@jfeala) for Outlier Bio (@outlierbio)

**exception BatchJobException**

Bases: `Exception`

---

```
class BatchClient(poll_time=10, **kwargs)
    Bases: object

    get_active_queue()
        Get name of first active job queue

    get_job_id_from_name(job_name)
        Retrieve the first job ID matching the given name

    get_job_status(job_id)
        Retrieve task statuses from ECS API

            Parameters (str) (job_id) – AWS Batch job uuid

            Returns one of {SUBMITTED|PENDING|RUNNABLE|STARTING|RUNNING|SUCCEEDED|FAILED}

    get_logs(log_stream_name, get_last=50)
        Retrieve log stream from CloudWatch

    submit_job(**kwargs)
        Wrap submit_job with useful defaults

    terminate_job(**kwargs)
        Wrap terminate_job

    hard_terminate(job_ids, reason, iattempt=0, **kwargs)
        Terminate all jobs with a hard(ish) exit via an Exception. The function will also wait for jobs to be explicitly terminated

    wait_on_job(job_id)
        Poll task status until STOPPED

    register_job_definition(json_fpath)
        Register a job definition with AWS Batch, using a JSON
```

## misctools

A collection of miscellaneous tools.

```
get_config(file_name, header, path='core/config/')
    Get the configuration from a file in the luigi config path directory, and convert the key-value pairs under the config header into a dict.
```

### Parameters

- **file\_name** (*str*) – The configuration file name.
- **header** (*str*) – The header key in the config file.

### Returns dict

```
get_paths_from_relative(relative=1)
```

A helper method for within `find_filepath_from_pathstub`. Prints all file and directory paths from a relative number of ‘backward steps’ from the current working directory.

```
find_filepath_from_pathstub(path_stub)
```

Find the full path of the ‘closest’ file (or directory) to the current working directory ending with `path_stub`. The *closest* file is determined by starting forwards of the current working directory. The algorithm is then repeated by moving the current working directory backwards, one step at a time until the file (or directory) is found. If the HOME directory is reached, the algorithm raises `FileNotFoundException`.

**Parameters path\_stub** (*str*) – The partial file (or directory) path stub to find.

**Returns** The full path to the partial file (or directory) path stub.

### mysqlDb

NOTE: overwhelmingly based on [this2](#), where the following documentation has been directly lifted. The main difference to the latter, is that `**cnx_kwargs` in the constructor can accept `port` as a key.

**class MySqlTarget** (*host, database, user, password, table, update\_id, \*\*cnx\_kwargs*)

Bases: luigi.target.Target

Target for a resource in MySql.

**marker\_table** = 'table\_updates'

**touch** (*connection=None*)

Mark this update as complete.

IMPORTANT, If the marker table doesn't exist, the connection transaction will be aborted and the connection reset. Then the marker table will be created.

**exists** (*connection=None*)

Returns True if the Target exists and False otherwise.

**connect** (*autocommit=False*)

**create\_marker\_table** ()

Create marker table if it doesn't exist.

Using a separate connection since the transaction might have to be reset.

### s3

A more recent implementation of AWS S3 support, stolen from: [https://gitlab.com/ced/s3\\_helpers/blob/master/luigi\\_s3\\_target.py](https://gitlab.com/ced/s3_helpers/blob/master/luigi_s3_target.py), but instead using modern boto3 commands.

**merge\_dicts** (\*dicts)

Merge dicts together, with later entries overriding earlier ones.

**parse\_s3\_path** (*path*)

For a given S3 path, return the bucket and key values

**class S3FS** (\*\*kwargs)

Bases: luigi.target.FileSystem

**exists** (*path*)

Return true if S3 key exists

**remove** (*path, recursive=True*)

Remove a file or directory from S3

**mkdir** (*path, parents=True, raise\_if\_exists=False*)

Create directory at location path

Creates the directory at path and implicitly create parent directories if they do not already exist.

#### Parameters

- **path** (*str*) – a path within the FileSystem to create as a directory.
- **parents** (*bool*) – Create parent directories when necessary. When parents=False and the parent directory doesn't exist, raise luigi.target.MissingParentDirectory

- **`raise_if_exists`** (`bool`) – raise `luigi.target.FileAlreadyExists` if the folder already exists.

**`isdir`** (`path`)

Return `True` if the location at `path` is a directory. If not, return `False`.

**Parameters** `path` (`str`) – a path within the `FileSystem` to check as a directory.

*Note:* This method is optional, not all `FileSystem` subclasses implements it.

**`listdir`** (`path`)

Return a list of files rooted in `path`.

This returns an iterable of the files rooted at `path`. This is intended to be a recursive listing.

**Parameters** `path` (`str`) – a path within the `FileSystem` to list.

*Note:* This method is optional, not all `FileSystem` subclasses implements it.

**`copy`** (`path, dest`)

Copy a file or a directory with contents. Currently, `LocalFileSystem` and `MockFileSystem` support only single file copying but `S3Client` copies either a file or a directory as required.

**`move`** (`path, dest`)

Move a file, as one would expect.

**`du`** (`path`)

**class** `S3Target` (`path, s3_args={}, **kwargs`)  
Bases: `luigi.target.FileSystemTarget`

**`fs = None`****`open`** (`mode='rb'`)

Open the `FileSystem` target.

This method returns a file-like object which can either be read from or written to depending on the specified mode.

**Parameters** `mode` (`str`) – the mode `r` opens the `FileSystemTarget` in read-only mode, whereas `w` will open the `FileSystemTarget` in write mode. Subclasses can implement additional options.

**class** `AtomicS3File` (`path, s3_obj, **kwargs`)  
Bases: `luigi.target.AtomicLocalFile`

**`move_to_final_destination()`**

## 2.2.6 Scripts

## 2.2.7 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.

## **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.

## **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.

## **2.2.8 Elasticsearch**

The following steps will take you through setting up elasticsearch on an EC2 instance.

Launch the EC2 instance and ssh in so the following can be installed:

### **docker**

```
sudo yum install docker -y
```

### **docker-compose**

```
curl -L https://github.com/docker/compose/releases/download/1.22.0/docker-compose-`uname -s` - `uname -m` -o /usr/local/bin/docker-compose
chmod +x /usr/local/bin/docker-compose
sudo ln -s /usr/local/bin/docker-compose /usr/bin/docker-compose
```

more info: <https://github.com/docker/compose/releases>

### **docker permissions**

```
sudo usermod -a -G docker $USER
```

more info: <https://techoverflow.net/2017/03/01/solving-docker-permission-denied-while-trying-to-connect-to-the-docker-daemon-socket/>

### **vm.max\_map\_count**

set permanantly in `/etc/sysctl.conf` by adding the following line: `vm.max_map_count=262144`

more info: <https://www.elastic.co/guide/en/elasticsearch/reference/current/docker.html>

## python 3.6

```
sudo yum install python36 -y
The machine now needs to be rebooted sudo reboot now
```

## Docker

**the docker-compose.yml needs to include ulimits settings::**

**ulimits:**

**memlock:** soft: -1 hard: -1

**nofile:** soft: 65536 hard: 65536

Recipes for http or https clusters can be found in: `nesta/core/scripts/elasticsearch`

There is also an EC2 AMI for a http node stored in the London region: `elasticsearch node London vx`

## Reindexing data from a remote cluster

- reindex permissions need to be set in the new cluster's `elasticsearch.yml`
- if the existing cluster is AWS hosted ES the ip address needs to be added to the security settings
- follow this guide: <https://www.elastic.co/guide/en/elasticsearch/reference/current/docs-reindex.html#reindex-from-remote>
- `index` and `query` do not need to be supplied
- if reindexing from AWS ES the port should be 443 for https. This is mandatory in the json sent to the reindexing api

## 2.2.9 Containerised Luigi

### Requirements

To containerise a pipeline a few steps are required:

- All imports must be absolute, ie `from nesta.` packages, core etc
- Once testing is complete the code should be committed and pushed to github, as this prevents the need to use local build options
- If building and running locally, Docker must be installed on the local machine and given enough RAM in the settings to run the pipeline
- Any required configuration files must be in `nesta.core.config` ie luigi and mysql config files, any API keys. **This directory is ignored but check before committing**

### Build

The build uses a multi-stage Dockerfile to reduce the size of the final image:

1. Code is git cloned and requirements are pip installed into a virtual environment
2. The environment is copied into the second image

From the root of the repository: `docker build -f docker/Dockerfile -t name:tag .`

Where `name` is the name of the created image and `tag` is the chosen tag. Eg `arxlive:latest`. This just makes the run step easier rather than using the generated id

The two stage build will normally just rebuild the second stage pulling in new code only. If a full rebuild is required, eg after `requirements.txt` has changed then include: `--no-cache`

Python version defaults to 3.6 but can be set during build by including the flag: `--build-arg python-version=3.7`

Tag defaults to `dev` but this can be overridden by including the flag: `--build-arg GIT_TAG=0.3` a branch name also works `--build-arg GIT_TAG=my-branch`

Work from a branch or locally while testing. Override the target branch from Github using the method above, or use the commented out code in the Dockerfile to switch to build from local files. **Don't commit this change!**

### Run

As only one pipeline runs in the container the `luigid` scheduler is not used.

There is a `docker-compose` file for arXlive which mounts your local `~/.aws` folder for AWS credentials as this outside docker's context:

```
docker-compose -f docker/docker-compose.yml run luigi --module module_path  
params
```

Where:

- `docker-compose.yml` is the docker-compose file containing the image: `image_name:tag` from the build
- `module_path` is the full python path to the module
- `params` are any other params to supply as per normal, ie `--date --production` etc

Eg: `docker-compose -f docker/docker-compose-arxlive-dev.yml run luigi --module nesta.core.routines.arxiv.arxiv_iterative_root_task RootTask --date 2019-04-16`

This could be adapted for each pipeline, or alternatively run with the volume specified with `-v`

```
docker run -v ~/.aws:/root/.aws:ro name:tag --module ...
```

Where `name` is the name of the created image and `tag` is the chosen tag. Eg `arxlive:latest` `--module ...` onwards contains the arguments you would pass to Luigi.

### Scheduling

1. Create an executable shell script in `nesta.core.scripts` to launch docker-compose with all the necessary parameters eg: `production`
2. Add a cron job to the shell script (there are some good online cron syntax checking tools, if needed)
3. Set the cron job to run every few minutes while testing and check the logs with `docker logs mycontainerhash --tail 50`. Obtain the hash using `docker ps`
4. It will run logged in as the user who set it up but there still may still be some permissions issues to deal with

## Currently scheduled

arXlive:

- A shell script to launch docker-compose for arXlive is set up to run in a cron job on user `russellwinch`
- This is scheduled for Sunday-Thursday at 0300 GMT. arXiv is updated on these days at 0200 GMT
- Logs are just stored in the container, use `docker logs container_name` to view

## Important points

- Keep any built images secure, they contain credentials
- You need to rebuild if code has changed
- As there is no central scheduler there is nothing stopping you from running the task more than once at the same time, by launching the container multiple times
- The graphical interface is not enabled without the scheduler

## Debugging

If necessary, it's possible to debug inside the container, but the endpoint needs to be overridden with bash:

```
docker run --entrypoint /bin/bash -itv ~/.aws:/root/.aws:ro image_name:tag
```

Where `image_name:tag` is the image from the build step. This includes the mounting of the `.aws` folder.

Almost nothing is installed (not even vi!!) other than Python so `apt-get update` and then `apt-get install` whatever you need

## Todo

A few things are sub-optimal:

- The container runs on the prod box rather than in the cloud in ECS
- Credentials are held in the container and local AWS config is required, this is the cause of the above point
- Due to the Nesta monorepo everything is pip installed, making a large container size with many unrequired packages. Pipeline specific requirements should be considered
- As logs are stored in the old containers they are kept until the next run where they are pruned and the logs are lost. Add a method of getting the logs to the host logger and record centrally
- In the arXlive pipeline there are at least 500 calls to the MAG API each run as the process tries to pick up new title matches on existing articles. As the API key only allows 10,000 calls per month this is currently OK with the schedule as it is but could possibly go over at some point



# CHAPTER 3

---

## AWS FAQ

---

### 3.1 Where is the data?

As a general rule-of-thumb, our data is always stored in the London region (`eu-west-2`), in either RDS (tier-0, MySQL) or Elasticsearch (tier-1). For the EURITO project we also use Neo4j (tier-1), and in the distant future we will use Neo4j for tier-2 (i.e. a knowledge graph).

### 3.2 Why don't you use Aurora rather than MySQL?

Aurora is definitely cheaper for stable production and business processes but not for research and development. You are charged for every byte of data you have *ever* consumed. This quickly spirals out-of-control for big-data development. Maybe one day we'll consider migrating back, once the situation stabilises.

### 3.3 Where are the production machines?

Production machines (EC2) run in Ohio (us-east-2).

### 3.4 Where is the latest config?

One of our priorities is to implement a decent config management system. Our latest read-only config for accessing the tier-0 (rawish SQL data) can be found [here](#), and a fairly up-to-date config directory (which you can paste into `nesta/core/config`) can be found [here](#). If you want to use exactly what Joel has been using, feel free to `sudo cp` from `/home/ec2-user/nesta/nesta/core/config`. For example, you can find the latest Elasticsearch indexes and endpoints [here](#).

## 3.5 Where do I start with Elasticsearch?

All Elasticsearch indexes (aka “databases” to the rest of the world), mappings (aka “schemas”) and whitelisting can be found [here](#).

I’d recommend using PostMan for spinning up and knocking down indexes. Practice this on a new cluster (which you can spin up from the above link), and then practice `PUT`, `POST` and `DELETE` requests to `PUT` an index (remember: “database”) with a mapping (“schema”), inserting a “row” with `POST` and then deleting the index with `DELETE`. You will quickly learn that it’s very easy to delete everything in Elasticsearch.

# CHAPTER 4

---

## Troubleshooting

---

### 4.1 How do I restart the apache server after downtime?

```
sudo service httpd restart
```

### 4.2 How do I restart the luigi server after downtime?

```
sudo su - luigi
source activate py36
luigid --background --pidfile /var/run/luigi/luigi.pid --logdir /var/log/luigi
```

### 4.3 How do I perform initial setup to ensure the batchables will run?

- AWS CLI needs to be installed and configured:

```
pip install awscli
aws configure
```

AWS Access Key ID and Secret Access Key are set up in IAM > Users > Security Credentials Default region name should be `eu-west-1` to enable the error emails to be sent In AWS SES the sender and receiver email addresses need to be verified

- The config files need to be accessible and the PATH and LUIGI\_CONFIG\_PATH need to be amended accordingly

### 4.4 How do I add a new user to the server?

- add the user with `useradd --create-home username`

- add sudo privileges following these instructions
- add to ec2 user group with `sudo usermod -a -G ec2-user username`
- set a temp password with `passwd username`
- their home directory will be `/home/username/`
- copy `.bashrc` to their home directory
- create folder `.ssh` in their home directory
- copy `.ssh/authorized_keys` to the same folder in their home directory (DONT MOVE IT!!)
- `cd` to their home directory and perform the below
- `chown` their copy of `.ssh/authorized_keys` to their username: `chown username .ssh/authorized_keys`
- clone the `nesta` repo
- copy `core/config` files
- set password to be changed next login `chage -d 0 username`
- share the temp password and core pem file

If necessary: `- sudo chmod g+w /var/tmp/batch`

# CHAPTER 5

---

## Packages

---

Nesta's collection of tools for meaty tasks. Any processes that go into production come here first, but there are other good reasons for code to end up here.



# CHAPTER 6

---

## Production

---

Nesta's production system is based on Luigi pipelines, and are designed to be entirely run on AWS via the batch service. The main Luigi server runs on a persistent EC2 instance. Beyond the well documented Luigi code, the main features of the nesta production system are:

- luigihacks.autobatch, which facilitates a managed Luigi.Task which is split, batched and combined in a single step. Currently only synchronous jobs are accepted. Asynchronous jobs (where downstream Luigi.Task jobs can be triggered) are a part of a longer term plan.
- scripts.nesta\_prepare\_batch which zips up the batchable with the specified environmental files and ships it to AWS S3.
- scripts.nesta\_docker\_build which builds a specified docker environment and ships it to AWS ECS.

### 6.1 License

MIT © 2018 Nesta



---

## Python Module Index

---

### C

core, 13  
core.batchables, 46  
core.batchables.batchgeocode.run, 51  
core.batchables.cordis.cordis\_api.run, 47  
core.batchables.crunchbase.crunchbase\_coffee.run, 47  
core.batchables.crunchbase.elasticsearch.run, 48  
core.batchables.examples.batch\_example.run, 46  
core.batchables.examples.template\_batchable.run, 46  
core.batchables.health\_data.nih\_abstract.run, 49  
core.batchables.health\_data.nih\_collect.run, 49  
core.batchables.health\_data.nih\_dedupe.run, 49  
core.batchables.health\_data.nih\_process.run, 49  
core.batchables.meetup.country\_groups.run, 50  
core.batchables.meetup.group\_details.run, 50  
core.batchables.meetup.groups\_members.run, 50  
core.batchables.meetup.members\_groups.run, 50  
core.batchables.meetup.topic\_tag\_elasticsearch.run, 50  
core.luighacks, 56  
core.luighacks.autobatch, 56  
core.luighacks.batchclient, 58  
core.luighacks.misctools, 59  
core.luighacks.mysql, 60  
core.luighacks.s3, 60  
core.orms.meetup\_orm, 52  
core.orms.nih\_orm, 53  
core.routines, 14  
core.routines.examples.batch\_example.batch\_example, 16  
core.routines.examples.db\_example.db\_example, 15  
core.routines.examples.s3\_example.s3\_example, 14  
nesta.core.batchables.arxiv.arxiv\_elasticsearch.run, 47  
nesta.core.batchables.eurito.arxiv\_eu.run, 48  
nesta.core.batchables.eurito.crunchbase\_eu.run, 48  
nesta.core.batchables.eurito.patstat\_eu.run, 48  
nesta.core.batchables.gtr.collect\_gtr.run, 48  
nesta.core.batchables.gtr.embed\_topics.run, 49  
nesta.core.batchables.nlp.corex\_topic\_model.run, 51  
nesta.core.batchables.nlp.ngrammer.run, 51  
nesta.core.batchables.nlp.tfidf.run, 51  
nesta.core.batchables.nlp.vectorizer.run, 51  
nesta.core.batchables.novelty.lolvelty.run, 18  
nesta.core.routines.arxiv.arxiv\_es\_task, 23  
nesta.core.routines.arxiv.arxiv\_es\_tokens, 24  
nesta.core.routines.arxiv.arxiv\_grid\_task, 21

```
nesta.core.routines.arxiv.arxiv_iterative_routine, 18
nesta.core.routines.arxiv.arxiv_lovelty, 24
nesta.core.routines.arxiv.arxiv_mag_sparqle_task, 20
nesta.core.routines.arxiv.arxiv_mag_task, 19
nesta.core.routines.arxiv.arxiv_root_task, 17
nesta.core.routines.arxiv.arxiv_topic_task, 22
nesta.core.routines.arxiv.deepchange_analyse, 25
nesta.core.routines.cordis.collect_cordis, 27
nesta.core.routines.crunchbase.crunchbase_collect, 33
nesta.core.routines.crunchbase.crunchbase_collect, 29
nesta.core.routines.crunchbase.crunchbase_collect, 31
nesta.core.routines.crunchbase.crunchbase_collect, 34
nesta.core.routines.crunchbase.crunchbase_collect, 28
nesta.core.routines.crunchbase.crunchbase_collect, 32
nesta.core.routines.crunchbase.crunchbase_non_org_collect_task, 29
nesta.core.routines.crunchbase.crunchbase_org_collect_task, 32
nesta.core.routines.crunchbase.crunchbase_parent_id_collect_task, 28
nesta.core.routines.crunchbase.crunchbase_root_task, 35
nesta.core.routines.eurito_es.es_root, 36
nesta.core.routines.eurito_es.preprocess_patstat, 37
nesta.core.routines.gtr.gtr_collect, 38
nesta.core.routines.gtr.gtr_geocode, 39
nesta.core.routines.gtr.gtr_root_task, 36
nesta.core.routines.health_data.nih_data.nih_abstracts_mesh_task, 41
nesta.core.routines.health_data.nih_data.nih_collect_task, 41
nesta.core.routines.health_data.nih_data.nih_dedupe_task, 41
nesta.core.routines.health_data.nih_data.nih_lovelty, 38
nesta.core.routines.health_data.nih_data.nih_process_task, 37
nesta.core.routines.health_data.nih_data.nih_root_task, 37
```

---

## Index

---

### A

abstract\_text (*Abstracts attribute*), 54  
Abstracts (*class in core.orms.nih\_orm*), 54  
AbstractsMeshTask (*class in nestar.core.routines.health\_data.nih\_data.nih\_abstracts\_mesh\_task*), 40  
activity (*Projects attribute*), 53  
administering\_ic (*Projects attribute*), 53  
affiliation (*Publications attribute*), 55  
AnalysisTask (*class in nestar.core.routines.arxiv.deepchange\_analysis\_task*), 25  
application\_id (*Abstracts attribute*), 54  
application\_id (*Projects attribute*), 53  
application\_type (*Projects attribute*), 53  
arra\_funded (*Projects attribute*), 53  
articles\_from\_date (*AnalysisTask attribute*), 26  
articles\_from\_date (*CollectNewTask attribute*), 18  
articles\_from\_date (*DateTask attribute*), 19  
articles\_from\_date (*GridTask attribute*), 22  
articles\_from\_date (*MagSparqlTask attribute*), 21  
articles\_from\_date (*QueryMagTask attribute*), 20  
articles\_from\_date (*RootTask attribute*), 17  
articles\_from\_date (*StandaloneAnalysisTask attribute*), 26  
ArxivElasticsearchTask (*class in nestar.core.routines.arxiv.arxiv\_lolvelty*), 24  
ArxivESTask (*class in nestar.core.routines.arxiv.arxiv\_es\_task*), 23  
ArxivESTokenTask (*class in nestar.core.routines.arxiv.arxiv\_es\_tokens*), 24  
ArxivLolveltyRootTask (*class in nestar.core.routines.arxiv.arxiv\_lolvelty*), 25  
ArxivTokenRootTask (*class in*

*nestar.core.routines.arxiv.arxiv\_es\_tokens*), 24  
assert\_iso2\_key () (*in module packages.meetup.country\_groups*), 4  
AtomicCSSEFile (*class in core.luighacks.s3*), 61  
author\_list (*Publications attribute*), 55  
author\_name (*Publications attribute*), 55  
AutoBatchTask (*class in core.luighacks.autobatch*), 56  
award\_notice\_date (*Projects attribute*), 53

### B

batch\_limits () (*ProcessTask method*), 39  
batchable (*AutoBatchTask attribute*), 57  
BatchClient (*class in core.luighacks.batchclient*), 58  
BatchJobException, 58  
bucket (*HealthLabelTask attribute*), 31  
budget\_end (*Projects attribute*), 53  
budget\_start (*Projects attribute*), 53

### C

category (*CountryGroupsTask attribute*), 42  
category (*GroupDetailsTask attribute*), 43  
category (*GroupsMembersTask attribute*), 43  
category (*MembersGroupsTask attribute*), 43  
category (*RootTask attribute*), 44  
category\_id (*Group attribute*), 53  
category\_name (*Group attribute*), 52  
category\_shortname (*Group attribute*), 52  
cfda\_code (*Projects attribute*), 53  
cherry\_picked (*AnalysisTask attribute*), 26  
cherry\_picked (*StandaloneAnalysisTask attribute*), 27  
cherry\_picked (*WriteTopicTask attribute*), 23  
chunker () (*in module nestar.core.batchables.nlp.tfidf.run*), 51  
chunks () (*in module nestar.core.routines.meetup.country\_extended\_groups.country\_ext* 42

chunksize (*PrepareArxivS3Data attribute*), 22  
 city (*Group attribute*), 52  
 classifier\_key (*HealthLabelTask attribute*), 31  
 clean\_abstract() (in module *core.batchables.health\_data.nih\_abstract\_mesh\_data*), 49  
 clean\_and\_tokenize() (in module *ages.nlp\_utils.preprocess*), 8  
 CollectNewTask (class in *nesta.core.routines.arxiv.arxiv\_collect\_iterative\_task*), 18  
 CollectTask (class in *nesta.core.routines.health\_data.nih\_data.nih\_collect\_task*), 38  
 combine() (*AbstractsMeshTask method*), 41  
 combine() (*AutoBatchTask method*), 57  
 combine() (*CollectTask method*), 38  
 combine() (*CordisCollectTask method*), 27  
 combine() (*CountryGroupsTask method*), 42  
 combine() (*CrunchbaseSql2EsTask method*), 34  
 combine() (*DedupeTask method*), 41  
 combine() (*FundingRoundGeocodeTask method*), 30  
 combine() (*GroupDetailsTask method*), 44  
 combine() (*GroupsMembersTask method*), 43  
 combine() (*GtrTask method*), 36  
 combine() (*InvestorGeocodeTask method*), 30  
 combine() (*IpoGeocodeTask method*), 30  
 combine() (*MembersGroupsTask method*), 43  
 combine() (*NonOrgCollectTask method*), 29  
 combine() (*OrgGeocodeTask method*), 30  
 combine() (*PeopleGeocodeTask method*), 31  
 combine() (*ProcessTask method*), 39  
 combine() (*SomeBatchTask method*), 16  
 command\_line() (in module *core.luigihacks.autobatch*), 56  
 connect() (*MySQLTarget method*), 60  
 copy() (*S3FS method*), 61  
 CordisCollectTask (class in *nesta.core.routines.cordis.collect\_cordis\_task*), 27  
 core (module), 13  
 core.batchables (module), 46  
 core.batchables.batchgeocode.run (module), 51  
 core.batchables.cordis.cordis\_api.run (module), 47  
 core.batchables.crunchbase.crunchbase\_collect (module), 47  
 core.batchables.crunchbase.crunchbase\_elasticsearch\_root\_task (module), 48  
 core.batchables.examples.batch\_example.run (module), 46  
 core.batchables.examples.template\_batchable.run (module), 46  
 core.batchables.health\_data.nih\_abstract\_mesh\_data (module), 49  
 core.batchables.health\_data.nih\_collect\_data.run (module), 49  
 core.batchables.health\_data.nih\_dedupe.run (module), 49  
 core.batchables.health\_data.nih\_process\_data.run (module), 49  
 core.batchables.meetup.country\_groups.run (module), 50  
 core.batchables.meetup.group\_details.run (module), 50  
 core.batchables.meetup.groups\_members.run (module), 50  
 core.batchables.meetup.members\_groups.run (module), 50  
 core.batchables.meetup.topic\_tag\_elasticsearch.run (module), 50  
 core.luigihacks (module), 56  
 core.luigihacks.autobatch (module), 56  
 core.luigihacks.batchclient (module), 58  
 core.luigihacks.misctools (module), 59  
 core.luigihacks.mysql (module), 60  
 core.luigihacks.s3 (module), 60  
 core.orms.meetup\_orm (module), 52  
 core.orms.nih\_orm (module), 53  
 core.routines (module), 14  
 core.routines.examples.batch\_example.batch\_example (module), 16  
 core.routines.examples.db\_example.db\_example (module), 15  
 core.routines.examples.s3\_example.s3\_example (module), 14  
 core\_categories (*MeetupHealthSql2EsTask attribute*), 45  
 core\_categories (*RootTask attribute*), 45  
 core\_categories (*TopicDiscoveryTask attribute*), 44  
 core\_project\_num (*Projects attribute*), 53  
 country (*Group attribute*), 52  
 country (*Publications attribute*), 55  
 country\_code (*MeetupCountryGroups attribute*), 4  
 country\_name (*Group attribute*), 53  
 CountryGroupsTask (class in *nesta.core.routines.meetup.country\_extended\_groups.country\_extended\_group*), 42  
 create\_marker\_table() (*MySQLTarget method*), 60  
 created (*Group attribute*), 52  
 CrunchbaseSql2EsTask (class in *nesta.core.routines.crunchbase.crunchbase\_elasticsearch\_root\_task*), 34  
 CrunchbaseSql2EsTask (class in *nesta.core.routines.crunchbase.crunchbase\_elasticsearch\_task*), 33

**D**

data\_path (*WriteTopicTask attribute*), 23  
 date (*AbstractsMeshTask attribute*), 40  
 date (*AnalysisTask attribute*), 26  
 date (*ArxivElasticsearchTask attribute*), 24  
 date (*ArxivESTask attribute*), 23  
 date (*ArxivLoveltyRootTask attribute*), 25  
 date (*ArxivTokenRootTask attribute*), 24  
 date (*CollectNewTask attribute*), 18  
 date (*CollectTask attribute*), 38  
 date (*CrunchbaseLoveltyRootTask attribute*), 34  
 date (*CrunchbaseSql2EsTask attribute*), 33  
 date (*DateTask attribute*), 19  
 date (*DedupeTask attribute*), 41  
 date (*DescriptionMeshTask attribute*), 33  
 date (*EsOnlyRootTask attribute*), 18  
 date (*FinalTask attribute*), 14  
 date (*FundingRoundGeocodeTask attribute*), 30  
 date (*GridTask attribute*), 21  
 date (*GtrGeocode attribute*), 37  
 date (*GtrOnlyRootTask attribute*), 37  
 date (*GtrTask attribute*), 36  
 date (*HealthLabelTask attribute*), 31  
 date (*InputData attribute*), 15  
 date (*InvestorGeocodeTask attribute*), 30  
 date (*IpoGeocodeTask attribute*), 30  
 date (*MagSparqlTask attribute*), 20  
 date (*MeetupLoveltyRootTask attribute*), 46  
 date (*NiHLoveltyRootTask attribute*), 42  
 date (*NonOrgCollectTask attribute*), 29  
 date (*OrgCollectTask attribute*), 29  
 date (*OrgGeocodeTask attribute*), 30  
 date (*ParentIdCollectTask attribute*), 32  
 date (*PatstatPreprocessRootTask attribute*), 35  
 date (*PeopleGeocodeTask attribute*), 31  
 date (*PreprocessPatstatTask attribute*), 35  
 date (*ProcessRootTask attribute*), 39  
 date (*ProcessTask attribute*), 38  
 date (*QueryMagTask attribute*), 20  
 date (*RootTask attribute*), 16, 17, 28, 35, 36, 38, 44, 45  
 date (*SomeBatchTask attribute*), 16  
 date (*SomeTask attribute*), 14, 15  
 date (*StandaloneAnalysisTask attribute*), 26  
 date (*WriteTopicTask attribute*), 23  
*DateTask* (class in *nesta.core.routines.arxiv.arxiv\_iterative*, 19  
 db\_conf\_env (*PrepareArxivS3Data attribute*), 22  
 db\_conf\_env (*WriteTopicTask attribute*), 23  
 db\_config (*InputData attribute*), 15  
 db\_config\_env (*AnalysisTask attribute*), 26  
 db\_config\_env (*CollectNewTask attribute*), 18  
 db\_config\_env (*CordisCollectTask attribute*), 27  
 db\_config\_env (*CrunchbaseSql2EsTask attribute*), 33

db\_config\_env (*DateTask attribute*), 19  
 db\_config\_env (*DescriptionMeshTask attribute*), 33  
 db\_config\_env (*GridTask attribute*), 21  
 db\_config\_env (*GtrGeocode attribute*), 37  
 db\_config\_env (*HealthLabelTask attribute*), 31  
 db\_config\_env (*MagSparqlTask attribute*), 21  
 db\_config\_env (*OrgCollectTask attribute*), 29  
 db\_config\_env (*ParentIdCollectTask attribute*), 32  
 db\_config\_env (*QueryMagTask attribute*), 20  
 db\_config\_env (*RootTask attribute*), 45  
 db\_config\_env (*StandaloneAnalysisTask attribute*), 26  
 db\_config\_env (*TopicDiscoveryTask attribute*), 44  
 db\_config\_path (*AbstractsMeshTask attribute*), 40  
 db\_config\_path (*AnalysisTask attribute*), 26  
 db\_config\_path (*CollectNewTask attribute*), 18  
 db\_config\_path (*CollectTask attribute*), 38  
 db\_config\_path (*CordisCollectTask attribute*), 27  
 db\_config\_path (*DateTask attribute*), 19  
 db\_config\_path (*DedupeTask attribute*), 41  
 db\_config\_path (*DescriptionMeshTask attribute*), 33  
 db\_config\_path (*EsOnlyRootTask attribute*), 18  
 db\_config\_path (*GridTask attribute*), 21  
 db\_config\_path (*MagSparqlTask attribute*), 21  
 db\_config\_path (*NonOrgCollectTask attribute*), 29  
 db\_config\_path (*ParentIdCollectTask attribute*), 32  
 db\_config\_path (*ProcessRootTask attribute*), 39  
 db\_config\_path (*ProcessTask attribute*), 38  
 db\_config\_path (*QueryMagTask attribute*), 20  
 db\_config\_path (*RootTask attribute*), 17, 38  
 db\_config\_path (*StandaloneAnalysisTask attribute*), 26  
 db\_config\_path (*WriteTopicTask attribute*), 23  
 debug (*RootTask attribute*), 17  
*DedupeTask* (class in *nesta.core.routines.health\_data.nih\_data.nih\_dedupe\_task*), 41  
 description (*Group attribute*), 53  
*DescriptionMeshTask* (class in *nesta.core.routines.crunchbase.crunchbase\_mesh\_task*), 32  
 direct\_cost\_amt (*Projects attribute*), 54  
 done\_check () (*AbstractsMeshTask method*), 40  
 done\_ids () (*ArxivElasticsearchTask method*), 24  
 done\_ids () (*ArxivESTokenTask method*), 24  
 drop\_and\_recreate (*AbstractsMeshTask attribute*), 40  
 drop\_and\_recreate (*ArxivElasticsearchTask attribute*), 24  
 drop\_and\_recreate (*ArxivESTask attribute*), 23  
 drop\_and\_recreate (*CrunchbaseSql2EsTask attribute*), 34  
 drop\_and\_recreate (*DedupeTask attribute*), 41

drop\_and\_recreate (*EsOnlyRootTask attribute*), 18  
 drop\_and\_recreate (*ProcessRootTask attribute*), 39  
 drop\_and\_recreate (*ProcessTask attribute*), 39  
 drop\_and\_recreate (*RootTask attribute*), 17, 28, 35, 38, 45  
 du () (*S3FS method*), 61

**E**

ed\_inst\_type (*Projects attribute*), 53  
 env\_files (*AutoBatchTask attribute*), 57  
 EsOnlyRootTask (class in *nesta.core.routines.arxiv.arxiv\_root\_task*), 17  
 execute () (*AutoBatchTask method*), 58  
 exists () (in module *nesta.core.routines.health\_data.nih\_data.nih\_collect\_task*), 38  
 exists () (*MySQLTarget method*), 60  
 exists () (*S3FS method*), 60  
 extract\_core\_orgs () (in module *core.batchables.cordis.cordis\_api.run*), 47  
 extract\_date () (in module *ages.format\_utils.datetools*), 10  
 extract\_year () (in module *ages.format\_utils.datetools*), 10  
 extract\_yearly\_funds () (in module *core.batchables.health\_data.nih\_dedupe.run*), 50

**F**

filter\_by\_idf () (in module *ages.nlp\_utils.preprocess*), 8  
 FinalTask (class in *core.routines.examples.s3\_example.s3\_example*), 14  
 find\_filepath\_from\_pathstub () (in module *core.luighacks.misctools*), 59  
 flatten\_data () (in module *ages.meetup.meetup\_utils*), 6  
 foa\_number (*Projects attribute*), 53  
 fs (*S3Target attribute*), 61  
 full\_project\_num (*Projects attribute*), 53  
 funding\_ics (*Projects attribute*), 53  
 funding\_mechanism (*Projects attribute*), 53  
 FundingRoundGeocodeTask (class in *nesta.core.routines.crunchbase.crunchbase\_geocode*), 30  
 fy (*Projects attribute*), 54

**G**

generate\_composite\_key () (in module *packages.geo\_utils.geocode*), 9  
 generate\_coords () (in module *ages.meetup.country\_groups*), 3

geocode () (in module *packages.geo\_utils.geocode*), 8  
 geocode\_batch\_dataframe () (in module *packages.geo\_utils.geocode*), 9  
 geocode\_dataframe () (in module *packages.geo\_utils.geocode*), 9  
 get\_abstract\_file\_keys () (*AbstractsMeshTask static method*), 40  
 get\_active\_queue () (*BatchClient method*), 59  
 get\_all\_members () (in module *packages.meetup.groups\_members*), 5  
 get\_api\_key () (in module *ages.meetup.meetup\_utils*), 6  
 get\_config () (in module *core.luighacks.misctools*), 59  
 get\_coordinate\_data () (in module *ages.meetup.country\_groups*), 4  
 get\_core\_topics () (in module *ages.meetup.meetup\_utils*), 7  
 get\_data\_urls () (in module *ages.health\_data.collect\_nih*), 7  
 get\_group\_details () (in module *ages.meetup.group\_details*), 6  
 get\_groups () (*MeetupCountryGroups method*), 4  
 get\_groups\_recursive () (*MeetupCountryGroups method*), 5  
 get\_job\_id\_from\_name () (*BatchClient method*), 59  
 get\_job\_status () (*BatchClient method*), 59  
 get\_logs () (*BatchClient method*), 59  
 get\_member\_details () (in module *ages.meetup.members\_groups*), 5  
 get\_member\_groups () (in module *ages.meetup.members\_groups*), 5  
 get\_members () (in module *ages.meetup.groups\_members*), 5  
 get\_members\_by\_percentile () (in module *packages.meetup.meetup\_utils*), 6  
 get\_paths\_from\_relative () (in module *core.luighacks.misctools*), 59  
 get\_value () (in module *core.batchables.health\_data.nih\_dedupe.run*), 50  
 grid\_task\_kwargs (*AnalysisTask attribute*), 26  
 grid\_task\_kwargs (*ArxivElasticsearchTask attribute*), 24  
 grid\_task\_kwargs (*ArxivESTask attribute*), 23  
 grid\_task\_kwargs (*PrepareArxivS3Data attribute*), 22  
 grid\_task\_kwargs (*StandaloneAnalysisTask attribute*), 27  
 grid\_task\_kwargs (*WriteTopicTask attribute*), 23  
 GridTask (class in *nesta.core.routines.arxiv.arxiv\_grid\_task*), 21  
 Group (class in *core.orms.meetup\_orm*), 52

group\_id (*GroupMember attribute*), 53  
 group\_urlname (*GroupMember attribute*), 53  
 GroupDetailsTask (class) in *nesta.core.routines.meetup.country\_extended\_groupextended\_geo(OrgCollectTask attribute)*, 29  
 43  
*nesta.core.routines.meetup.country\_extended\_groupextended\_geo(OrgCollectTask attribute)*, 29  
 43  
 GroupMember (*class in core.orms.meetup\_orm*), 53  
 groups (*MeetupCountryGroups attribute*), 4  
 GroupsMembersTask (class) in *nesta.core.routines.meetup.country\_extended\_groupextended\_geo(PeopleGeocodeTask attribute)*, 31  
 42  
 GtrGeocode (class) in *nesta.core.routines.gtr.gtr\_geocode*, 37  
 GtrOnlyRootTask (class) in *nesta.core.routines.gtr.gtr\_collect*, 36  
 GtrTask (*class in nesta.core.routines.gtr.gtr\_collect*), 36

**H**

hard\_terminate() (*BatchClient method*), 59  
 HealthLabelTask (class) in *nesta.core.routines.crunchbase.crunchbase\_healthlabeltask*, 31  
 hierarchy\_field() (in module *nesta.core.batchables.arxiv.arxiv\_elasticsearch.run*), 47

**I**

ic\_name (*Projects attribute*), 54  
 id (*Group attribute*), 52  
 in\_db\_config (*SomeTask attribute*), 15  
 index (*CrunchbaseLoveltyRootTask attribute*), 34  
 index (*MeetupLoveltyRootTask attribute*), 46  
 index (*NiHLovelyRootTask attribute*), 42  
 indirect\_cost\_amt (*Projects attribute*), 54  
 InputData (class) in *core.routines.examples.db\_example.db\_example*, 15  
 InputData (class) in *core.routines.examples.s3\_example.s3\_example*, 14  
 insert\_batch\_size (*AnalysisTask attribute*), 26  
 insert\_batch\_size (*CollectNewTask attribute*), 18  
 insert\_batch\_size (*CrunchbaseSql2EsTask attribute*), 33  
 insert\_batch\_size (*DateTask attribute*), 19  
 insert\_batch\_size (*DescriptionMeshTask attribute*), 33  
 insert\_batch\_size (*FundingRoundGeocodeTask attribute*), 30  
 insert\_batch\_size (*GridTask attribute*), 21  
 insert\_batch\_size (*HealthLabelTask attribute*), 31  
 insert\_batch\_size (*InvestorGeocodeTask attribute*), 30  
 insert\_batch\_size (*IpoGeocodeTask attribute*), 30  
 insert\_batch\_size (*MagSparqlTask attribute*), 21  
 insert\_batch\_size (*NonOrgCollectTask attribute*), 29  
 insert\_batch\_size (*OrgGeocodeTask attribute*), 30  
 insert\_batch\_size (*ParentIdCollectTask attribute*), 32  
 insert\_batch\_size (*nesta.core.routines.meetup.country\_extended\_groupextended\_geo(PeopleGeocodeTask attribute)*), 31  
 insert\_batch\_size (*QueryMagTask attribute*), 20  
 insert\_batch\_size (*RootTask attribute*), 17, 28  
 insert\_batch\_size (*StandaloneAnalysisTask attribute*), 26  
 insert\_batch\_size (*WriteTopicTask attribute*), 23  
 intermediate\_bucket (*CordisCollectTask attribute*), 27  
 intermediate\_bucket (*CrunchbaseSql2EsTask attribute*), 34  
 intermediate\_bucket (*DedupeTask attribute*), 41  
 IpoGeocodeTask (class) in *nesta.core.routines.crunchbase.crunchbase\_geocode\_task*, 30  
 isdir() (*S3FS method*), 61  
 iso2 (*CountryGroupsTask attribute*), 42  
 iso2 (*GroupDetailsTask attribute*), 43  
 iso2 (*GroupsMembersTask attribute*), 43  
 iso2 (*MembersGroupsTask attribute*), 43  
 iso2 (*RootTask attribute*), 44  
 issn (*Publications attribute*), 55  
 iterrows() (in module *nesta.health\_data.collect\_nih*), 7

**J**

job\_def (*AutoBatchTask attribute*), 57  
 job\_name (*AutoBatchTask attribute*), 57  
 job\_queue (*AutoBatchTask attribute*), 57  
 journal\_issue (*Publications attribute*), 55  
 journal\_title (*Publications attribute*), 55  
 journal\_title\_abbr (*Publications attribute*), 55  
 journal\_volume (*Publications attribute*), 55

**K**

kwarg\_maker() (in module *nesta.core.routines.eurito\_es.es\_root*), 35

**L**

lang (*Publications attribute*), 55  
 lat (*Group attribute*), 53  
 listdir() (*S3FS method*), 61

`load_transformer()` (in module `packag.es.decorators.schema_transform`), 10  
`lon` (*Group attribute*), 53

**M**

`mag_config_path` (*AnalysisTask attribute*), 26  
`mag_config_path` (*GridTask attribute*), 21  
`mag_config_path` (*MagSparqlTask attribute*), 21  
`mag_config_path` (*QueryMagTask attribute*), 20  
`mag_config_path` (*StandaloneAnalysisTask attribute*), 26  
`MagSparqlTask` (class) in `nesta.core.routines.arxiv.arxiv_mag_sparql_task`, 20  
`marker_table` (*MySqlTarget attribute*), 60  
`max_age` (*SomeTask attribute*), 15  
`max_live_jobs` (*AutoBatchTask attribute*), 57  
`max_runs` (*AutoBatchTask attribute*), 57  
`MeetupCountryGroups` (class) in `packag.es.meetup.country_groups`, 4  
`MeetupHealthSql2EsTask` (class) in `nesta.core.routines.meetup.health_tagging.health_meetup_es`, 45  
`MeetupLoveltyRootTask` (class) in `nesta.core.routines.meetup.health_tagging.meetup_lovelty`, 46  
`member_id` (*GroupMember attribute*), 53  
`members` (*Group attribute*), 53  
`members_perc` (*MeetupHealthSql2EsTask attribute*), 45  
`members_perc` (*RootTask attribute*), 45  
`members_perc` (*TopicDiscoveryTask attribute*), 44  
`MembersGroupsTask` (class) in `nesta.core.routines.meetup.country_extended_groups`, 43  
`memory` (*AutoBatchTask attribute*), 57  
`merge_dicts()` (in module `core.luighacks.s3`), 60  
`merge_lists()` (in module `nesta.core.batchables.nlp.vectorizer.run`), 52  
`metadata()` (in module `nesta.core.batchables.eurito.patstat_eu.run`), 48  
`mkdir()` (*S3FS method*), 60  
`move()` (*S3FS method*), 61  
`move_to_final_destination()` (*AtomicS3File method*), 61  
`MySqlTarget` (class in `core.luighacks.mysql_db`), 60

**N**

`name` (*Group attribute*), 53  
`nesta.core.batchables.arxiv.arxiv_elasticse` (module), 47

`nesta.core.batchables.eurito.arxiv_eu.run` (module), 48  
`nesta.core.batchables.eurito.crunchbase_eu.run` (module), 48  
`nesta.core.batchables.eurito.patstat_eu.run` (module), 48  
`nesta.core.batchables.gtr.collect_gtr.run` (module), 48  
`nesta.core.batchables.gtr.embed_topics.run` (module), 49  
`nesta.core.batchables.nlp.corex_topic_model.run` (module), 51  
`nesta.core.batchables.nlp.ngrammer.run` (module), 51  
`nesta.core.batchables.nlp.tfidf.run` (module), 51  
`nesta.core.batchables.nlp.vectorizer.run` (module), 51  
`nesta.core.batchables.novelty.lovelty.run` (module), 52  
`nesta.core.routines.arxiv.arxiv_collect_iterative_t` (module), 18  
`nesta.core.routines.arxiv.arxiv_es_task` (module), 23  
`nesta.core.routines.arxiv.arxiv_es_tokens` (module), 24  
`nesta.core.routines.arxiv.arxiv_grid_task` (module), 21  
`nesta.core.routines.arxiv.arxiv_iterative_date_task` (module), 18  
`nesta.core.routines.arxiv.arxiv_lovelty` (module), 24  
`nesta.core.routines.arxiv.arxiv_mag_sparql_task` (module), 19  
`nesta.core.routines.arxiv.arxiv_root_task` (module), 17  
`nesta.core.routines.arxiv.arxiv_topic_tasks` (module), 22  
`nesta.core.routines.arxiv.deepchange_analysis_task` (module), 25  
`nesta.core.routines.cordis.collect_cordis_task` (module), 27  
`nesta.core.routines.crunchbase.crunchbase_elasticse` (module), 33  
`nesta.core.routines.crunchbase.crunchbase_geocode_t` (module), 29  
`nesta.core.routines.crunchbase.crunchbase_health_l` (module), 31  
`nesta.core.routines.crunchbase.crunchbase_lovelty` (module), 34  
`nesta.core.routines.crunchbase.crunchbase_mesh_task` (module), 32

nesta.core.routines.crunchbase.crunchbase\_~~org~~  
~~on~~~~country~~(*Projects attribute*), 54  
     (*module*), 29  
 nesta.core.routines.crunchbase.crunchbase\_~~org~~  
~~register~~(*Projects attribute*), 54  
     (*module*), 28  
 nesta.core.routines.crunchbase.crunchbase\_~~org~~  
~~afépt~~(*Projects attribute*), 54  
     (*module*), 32  
 nesta.core.routines.crunchbase.crunchbase\_~~org~~  
~~oname~~(*Projects attribute*), 54  
     (*module*), 28  
 nesta.core.routines.eurito\_es.es\_root  
     (*module*), 35  
 nesta.core.routines.eurito\_es.preprocess\_patst  
     (*nesta.core.routines.crunchbase.crunchbase\_org\_collect\_task*),  
     28  
 nesta.core.routines.gtr.gtr\_collect  
     (*module*), 36  
 nesta.core.routines.gtr.gtr\_geocode  
     (*module*), 37  
 nesta.core.routines.gtr.gtr\_root\_task  
     (*module*), 36  
 nesta.core.routines.health\_data.nih\_dataon  
~~ipu~~  
~~qual~~(*CollectNewTask method*), 18  
     (*module*), 39  
 nesta.core.routines.health\_data.nih\_dataon  
~~ipu~~  
~~onl~~(*CountryCollectTask method*), 27  
     (*module*), 38  
 nesta.core.routines.health\_data.nih\_dataon  
~~ipu~~  
~~onl~~(*DatabaseSql2EsTask method*), 34  
     (*module*), 41  
 nesta.core.routines.health\_data.nih\_dataon  
~~ipu~~  
~~onl~~(*DescriptionMeshTask method*), 33  
     (*module*), 41  
 nesta.core.routines.health\_data.nih\_dataon  
~~ipu~~  
~~onl~~(*FindingRoundGeocodeTask method*), 30  
     (*module*), 38  
 nesta.core.routines.health\_data.nih\_dataon  
~~ipu~~  
~~onl~~(*GroupDetailsTask method*), 43  
     (*module*), 37  
 nesta.core.routines.meetup.country\_extended  
~~ogt~~  
~~on~~(*GtGeocode method*), 37  
     (*module*), 42  
 nesta.core.routines.meetup.health\_taggingu  
~~heal~~  
~~th~~(*HealthLabelTask method*), 32  
     (*module*), 45  
 nesta.core.routines.meetup.health\_taggingu  
~~heal~~  
~~th~~(*IpoGeocodeTask method*), 30  
     (*module*), 45  
 nesta.core.routines.meetup.health\_taggingu  
~~heal~~  
~~th~~(*MagSpqrqlTask method*), 21  
     (*module*), 44  
 nih\_spending\_cats (*Projects attribute*), 54  
 NIHLovertyRootTask  
     (*class*)  
     (*nesta.core.routines.health\_data.nih\_data.nih\_loverty*), 41  
 NoGroupFound, 6  
 NoMemberFound, 5  
 NonOrgCollectTask  
     (*class*)  
     (*nesta.core.routines.crunchbase.crunchbase\_non\_org\_collect\_task*), 29  
  
**O**  
 open () (*S3Target method*), 61  
 optional ()  
     (*in module*)  
     (*nesta.core.batchables.nlp.vectorizer.run*), 52  
 org\_city (*Projects attribute*), 54

**P**

packages (*module*), 3  
 packages.decorators.ratelimit (*module*), 10  
 packages.decorators.schema\_transform (*module*), 10  
 packages.format\_utils.datetools (*module*), 9  
 packages.geo\_utils.geocode (*module*), 8  
 packages.health\_data.collect\_nih (*module*), 7  
 packages.health\_data.process\_nih (*module*), 7  
 packages.meetup.country\_groups (*module*), 3  
 packages.meetup.group\_details (*module*), 6  
 packages.meetup.groups\_members (*module*), 5  
 packages.meetup.meetup\_utils (*module*), 6  
 packages.meetup.members\_groups (*module*), 5  
 packages.nlp\_utils.preprocess (*module*), 8  
 page\_number (*Publications attribute*), 55  
 page\_size (*GtrGeocode attribute*), 37  
 page\_size (*GtrOnlyRootTask attribute*), 37  
 page\_size (*GtrTask attribute*), 36  
 page\_size (*RootTask attribute*), 36  
 ParentIdCollectTask (*class* in *nesta.core.routines.crunchbase.crunchbase\_parent\_id\_collect\_task*), 32  
 parse\_s3\_path () (*in module core.luigihacks.s3*), 60  
 patent\_id (*Patents attribute*), 55  
 patent\_org\_name (*Patents attribute*), 55  
 patent\_title (*Patents attribute*), 55  
 Patents (*class in core.orms.nih\_orm*), 55  
 PatstatPreprocessRootTask (*class* in *nesta.core.routines.eurito\_es.preprocess\_patstat*), 35  
 PeopleGeocodeTask (*class* in *nesta.core.routines.crunchbase.crunchbase\_geocode\_task*), 31  
 phr (*Projects attribute*), 54  
 pi\_ids (*Projects attribute*), 54  
 pi\_names (*Projects attribute*), 54  
 pmc\_id (*Publications attribute*), 55  
 pmid (*Publications attribute*), 55  
 poll\_time (*AutoBatchTask attribute*), 57  
 prepare () (*AbstractsMeshTask method*), 40  
 prepare () (*AutoBatchTask method*), 57  
 prepare () (*CollectTask method*), 38  
 prepare () (*CordisCollectTask method*), 27  
 prepare () (*CountryGroupsTask method*), 42  
 prepare () (*CrunchbaseSql2EsTask method*), 34  
 prepare () (*DedupeTask method*), 41  
 prepare () (*GroupDetailsTask method*), 44  
 prepare () (*GroupsMembersTask method*), 43  
 prepare () (*GtrTask method*), 36  
 prepare () (*MembersGroupsTask method*), 43  
 prepare () (*NonOrgCollectTask method*), 29  
 prepare () (*ProcessTask method*), 39  
 prepare () (*SomeBatchTask method*), 16  
 prepare\_data () (*in module core.batchables.cordis.cordis\_api.run*), 47  
 PrepareArxivS3Data (*class* in *nesta.core.routines.arxiv.arxiv\_topic\_tasks*), 22  
 PreprocessPatstatTask (*class* in *nesta.core.routines.eurito\_es.preprocess\_patstat*), 35  
 process\_batch\_size (*CordisCollectTask attribute*), 27  
 process\_batch\_size (*CrunchbaseSql2EsTask attribute*), 33  
 process\_batch\_size (*DedupeTask attribute*), 41  
 process\_batch\_size (*RootTask attribute*), 35  
 ProcessRootTask (*class* in *nesta.core.routines.health\_data.nih\_data.nih\_process\_task*), 39  
 ProcessTask (*class* in *nesta.core.routines.health\_data.nih\_data.nih\_process\_task*), 38  
 production (*ArxivLoveltyRootTask attribute*), 25  
 production (*ArxivTokenRootTask attribute*), 24  
 production (*CrunchbaseLoveltyRootTask attribute*), 34  
 production (*EsOnlyRootTask attribute*), 18  
 production (*GtrOnlyRootTask attribute*), 37  
 production (*MeetupLoveltyRootTask attribute*), 46  
 production (*NiHLoveltyRootTask attribute*), 41  
 production (*PatstatPreprocessRootTask attribute*), 35  
 production (*ProcessRootTask attribute*), 39  
 production (*RootTask attribute*), 17, 27, 28, 35, 36, 38, 44, 45  
 production (*StandaloneAnalysisTask attribute*), 26  
 program\_officer\_name (*Projects attribute*), 54  
 project\_end (*Projects attribute*), 54  
 project\_id (*Patents attribute*), 55  
 project\_start (*Projects attribute*), 54  
 project\_terms (*Projects attribute*), 54  
 project\_title (*Projects attribute*), 54  
 Projects (*class in core.orms.nih\_orm*), 53  
 pub\_date (*Publications attribute*), 55  
 pub\_title (*Publications attribute*), 55  
 pub\_year (*Publications attribute*), 55  
 Publications (*class in core.orms.nih\_orm*), 54

**Q**

QueryMagTask (*class* in *nesta.core.routines.arxiv.arxiv\_mag\_task*), 19

**R**

ratelimit () (*in module pack-*)

ages.decorators.ratelimit), 10  
`raw_data_path` (*AnalysisTask* attribute), 26  
`raw_data_path` (*StandaloneAnalysisTask* attribute), 26  
`raw_s3_path_prefix` (*WriteTopicTask* attribute), 23  
`region_name` (*AutoBatchTask* attribute), 57  
`register_job_definition()` (*BatchClient* method), 59  
`remove()` (*S3FS* method), 60  
`requires()` (*AbstractsMeshTask* method), 40  
`requires()` (*AnalysisTask* method), 26  
`requires()` (*ArxivElasticsearchTask* method), 24  
`requires()` (*ArxivESTask* method), 23  
`requires()` (*ArxivLoveltyRootTask* method), 25  
`requires()` (*ArxivTokenRootTask* method), 24  
`requires()` (*CrunchbaseLoveltyRootTask* method), 34  
`requires()` (*CrunchbaseSql2EsTask* method), 34  
`requires()` (*DateTask* method), 19  
`requires()` (*DedupeTask* method), 41  
`requires()` (*DescriptionMeshTask* method), 33  
`requires()` (*EsOnlyRootTask* method), 18  
`requires()` (*FinalTask* method), 14  
`requires()` (*FundingRoundGeocodeTask* method), 30  
`requires()` (*GridTask* method), 22  
`requires()` (*GroupDetailsTask* method), 43  
`requires()` (*GroupsMembersTask* method), 43  
`requires()` (*GtrGeocode* method), 37  
`requires()` (*GtrOnlyRootTask* method), 37  
`requires()` (*HealthLabelTask* method), 31  
`requires()` (*InvestorGeocodeTask* method), 30  
`requires()` (*IpoGeocodeTask* method), 30  
`requires()` (*MagSparqlTask* method), 21  
`requires()` (*MeetupHealthSql2EsTask* method), 45  
`requires()` (*MeetupLoveltyRootTask* method), 46  
`requires()` (*MembersGroupsTask* method), 43  
`requires()` (*NiHLoveltyRootTask* method), 42  
`requires()` (*NonOrgCollectTask* method), 29  
`requires()` (*OrgGeocodeTask* method), 30  
`requires()` (*ParentIdCollectTask* method), 32  
`requires()` (*PatstatPreprocessRootTask* method), 35  
`requires()` (*PeopleGeocodeTask* method), 31  
`requires()` (*PrepareArxivS3Data* method), 22  
`requires()` (*ProcessRootTask* method), 39  
`requires()` (*ProcessTask* method), 39  
`requires()` (*QueryMagTask* method), 20  
`requires()` (*RootTask* method), 16, 17, 28, 35, 36, 38, 44, 45  
`requires()` (*SomeBatchTask* method), 16  
`requires()` (*SomeTask* method), 14, 15  
`requires()` (*StandaloneAnalysisTask* method), 27  
`requires()` (*WriteTopicTask* method), 23  
`retry_if_not_value_error()` (*in module* `packages.geo_utils.geocode`), 9  
`RootTask` (*class in core.routines.examples.batch\_example.batch\_example*), 16  
`RootTask` (*class in core.routines.examples.db\_example.db\_example*), 15  
`RootTask` (*class in nesta.core.routines.arxiv.arxiv\_root\_task*), 17  
`RootTask` (*class in nesta.core.routines.cordis.collect\_cordis\_task*), 27  
`RootTask` (*class in nesta.core.routines.crunchbase.crunchbase\_root\_task*), 28  
`RootTask` (*class in nesta.core.routines.eurito\_es.es\_root*), 35  
`RootTask` (*class in nesta.core.routines.gtr.gtr\_root\_task*), 36  
`RootTask` (*class in nesta.core.routines.health\_data.nih\_data.nih\_root\_task*), 37  
`RootTask` (*class in nesta.core.routines.meetup.country\_extended\_groups*), 44  
`RootTask` (*class in nesta.core.routines.meetup.health\_tagging.health\_meetup*), 45  
`routine_id` (*CordisCollectTask* attribute), 27  
`routine_id` (*DedupeTask* attribute), 41  
`routine_id` (*TopicDiscoveryTask* attribute), 44  
`run()` (*AnalysisTask* method), 26  
`run()` (*AutoBatchTask* method), 57  
`run()` (*CollectNewTask* method), 18  
`run()` (*DescriptionMeshTask* method), 33  
`run()` (*FinalTask* method), 15  
`run()` (*GridTask* method), 22  
`run()` (*GtrGeocode* method), 37  
`run()` (*HealthLabelTask* method), 32  
`run()` (*in module* `core.batchables.batchgeocode.run`), 51  
`run()` (*in module* `core.batchables.cordis.cordis_api.run`), 47  
`run()` (*in module* `core.batchables.crunchbase.crunchbase_collect.run`), 48  
`run()` (*in module* `core.batchables.crunchbase.crunchbase_elasticsearch.run`), 48  
`run()` (*in module* `core.batchables.examples.batch_example.run`), 46  
`run()` (*in module* `core.batchables.examples.template_batchable.run`), 47  
`run()` (*in module* `core.batchables.health_data.nih_abstract_mesh_data.run`), 49  
`run()` (*in module* `core.batchables.health_data.nih_collect_data.run`), 49  
`run()` (*in module* `core.batchables.health_data.nih_dedupe.run`), 50  
`run()` (*in module* `core.batchables.health_data.nih_process_data.run`), 49  
`run()` (*in module* `core.batchables.meetup.country_groups.run`), 50  
`run()` (*in module* `core.batchables.meetup.group_details.run`),

```

    50
run () (in module core.batchables.meetup.groups_members.run),   schema_transformer () (in module pack-
    50                                ages.decorators.schema_transform), 10
run () (in module core.batchables.meetup.members_groups.run),   select_text () (in module
    50                                nesto.core.batchables.eurito.patstat_eu.run),
    48
run () (in module core.batchables.meetup.topic_tag_elasticsearch.run), number (Projects attribute), 54
    51
run () (in module nesto.core.batchables.arxiv.arxiv_elasticsearch.run), SomeBatchTask (class
    47                                in
run () (in module nesto.core.routines.examples.batch_example.batch_example),
    16
run () (in module nesto.core.batchables.eurito.arxiv_eu.run), SomeInitialTask (class
    48                                in
run () (in module nesto.core.batchables.eurito.crunchbase_eu.run), 16
    48
run () (in module nesto.core.batchables.eurito.patstat_eu.run), SomeTask (class in core.routines.examples.db_example.db_example),
    15
run () (in module nesto.core.batchables.gtr.collect_gtr.run), SomeTask (class in core.routines.examples.s3_example.s3_example),
    48
run () (in module nesto.core.batchables.nlp.corex_topic_model.run), split_links () (in
    49                                module
run () (in module nesto.core.batchables.gtr.embed_topics.run), core.batchables.cordis.cordis_api.run), 47
    49
run () (in module nesto.core.routines.health_data.nih_data.nih_abstracts_mesh_task),
    51
run () (in module nesto.core.routines.health_data.nih_data.nih_abstracts_mesh_task),
    40
run () (in module nesto.core.routines.health_data.nih_data.nih_abstracts_mesh_task),
    51
run () (in module nesto.core.routines.health_data.nih_data.nih_abstracts_mesh_task),
    51
run () (in module nesto.core.routines.health_data.nih_data.nih_abstracts_mesh_task),
    25
run () (in module nesto.core.routines.health_data.nih_data.nih_abstracts_mesh_task),
    51
run () (in module nesto.core.routines.health_data.nih_data.nih_abstracts_mesh_task),
    26
run () (in module nesto.core.routines.health_data.nih_data.nih_abstracts_mesh_task),
    52
run () (in module nesto.core.routines.health_data.nih_data.nih_abstracts_mesh_task),
    52
run () (InputData method), 15
run () (MagSparqlTask method), 21
run () (OrgCollectTask method), 29
run () (ParentIdCollectTask method), 32
run () (PrepareArxivS3Data method), 23
run () (PreprocessPatstatTask method), 35
run () (QueryMagTask method), 20
run () (RootTask method), 17
run () (SomeTask method), 14, 15
run () (TopicDiscoveryTask method), 45
run () (WriteTopicTask method), 23

S
s3_path_out (PrepareArxivS3Data attribute), 22
s3_path_prefix (AnalysisTask attribute), 26
s3_path_prefix (StandaloneAnalysisTask attribute),
    26
s3_path_prefix (WriteTopicTask attribute), 23
S3FS (class in core.luighacks.s3), 60
S3Target (class in core.luighacks.s3), 61
save_sample () (in module pack-
    ages.meetup.meetup_utils), 6
schema_transform () (in module pack-
    ages.decorators.schema_transform), 10
    50
    50
    48
    51
    16
    15
    14
    49
    49
    40
    59
    54
    54
    40
    57
    54
    54
    54
    52
    52
    52
    59
    26
    57
    18
    19
    33
    21
    37
    31
    20
    29
    32

```

**T**

```

term_counts () (in module
nesto.core.routines.health_data.nih_data.nih_abstracts_mesh_task),
    52
terminate_job () (BatchClient method), 59
test (AnalysisTask attribute), 26
test (AutoBatchTask attribute), 57
test (CollectNewTask attribute), 18
test (DateTask attribute), 19
test (DescriptionMeshTask attribute), 33
test (GridTask attribute), 21
test (GtrGeocode attribute), 37
test (HealthLabelTask attribute), 31
test (MagSparqlTask attribute), 20
test (OrgCollectTask attribute), 29
test (ParentIdCollectTask attribute), 32

```

test (*PrepareArxivS3Data attribute*), 22  
test (*PreprocessPatstatTask attribute*), 35  
test (*QueryMagTask attribute*), 20  
test (*StandaloneAnalysisTask attribute*), 26  
test (*TopicDiscoveryTask attribute*), 44  
test (*WriteTopicTask attribute*), 23  
timeout (*AutoBatchTask attribute*), 57  
timestamp (*Group attribute*), 53  
tokenize\_document () (in module *pack-ages.nlp\_utils.preprocess*), 8  
topic\_perc (*MeetupHealthSql2EsTask attribute*), 45  
topic\_perc (*RootTask attribute*), 45  
topic\_perc (*TopicDiscoveryTask attribute*), 44  
TopicDiscoveryTask (class in *nesta.core.routines.meetup.health\_tagging.topic\_discovery\_task*), 44  
topics (*Group attribute*), 53  
total\_cost (*Projects attribute*), 54  
total\_cost\_sub\_project (*Projects attribute*), 54  
touch () (*MySqlTarget method*), 60

## U

urlname (*Group attribute*), 52

## V

vcpus (*AutoBatchTask attribute*), 57  
vectoriser\_key (*HealthLabelTask attribute*), 31

## W

wait\_on\_job () (*BatchClient method*), 59  
worker\_timeout (*AutoBatchTask attribute*), 57  
write\_to\_s3 () (*PrepareArxivS3Data method*), 22  
WriteTopicTask (class in *nesta.core.routines.arxiv.arxiv\_topic\_tasks*), 23