
BioThings SDK Documentation

Release 0.1

BioThings team

Dec 11, 2017

1	Introduction	1
1.1	What's BioThings?	1
1.2	BioThings SDK	1
1.3	BioThings API	1
2	Installing BioThings SDK	3
2.1	Single Data Source, No Source Updating Tutorial	3
2.1.1	Prerequisites	3
2.1.2	PharmGKB Gene	4
2.1.3	Generate BioThings API	5
2.2	Multiple Data Sources, Automated Source Updating Tutorial	7
2.2.1	Prerequisites	7
2.2.2	Configuration file	7
2.2.3	hub.py	8
2.2.4	Dumpers	9
2.2.5	Uploaders	14
2.2.6	Mergers	20
2.2.7	Indexers	29
2.3	BioThings Standalone instances	29
2.3.1	Quick Links	30
2.3.2	Prerequisites	30
2.3.3	What you'll learn	31
2.3.4	Data found in standalone instances	31
2.3.5	Downloading and running a standalone instance	31
2.3.6	Updating data using Biothings hub	33
2.3.7	BioThings API with multiple indices	38
2.3.8	Troubleshooting	38
2.4	Hub component	39
2.4.1	dumper	41
2.4.2	uploader	41
2.4.3	builder	41
2.4.4	indexer	41
2.4.5	differ	41
2.4.6	syncer	41
2.5	Web component	41
2.5.1	Server boot script	42

2.5.2	Settings	42
2.5.3	Handlers	43
2.5.4	Elasticsearch Query Builder	45
2.5.5	Elasticsearch Query	47
2.5.6	Elasticsearch Result Transformer	47

Python Module Index	49
----------------------------	-----------

1.1 What's BioThings?

We use “**BioThings**” to refer to objects of any biomedical entity-type represented in the biological knowledge space, such as genes, genetic variants, drugs, chemicals, diseases, etc.

1.2 BioThings SDK

SDK represents “Software Development Kit”. BioThings SDK provides a [Python-based](#) toolkit to build high-performance data APIs (or web services) from a single data source or multiple data sources. It has the particular focus on building data APIs for biomedical-related entities, a.k.a “*BioThings*”, though it's not necessarily limited to the biomedical scope. For any given “*BioThings*” type, BioThings SDK helps developers to aggregate annotations from multiple data sources, and expose them as a clean and high-performance web API.

The BioThings SDK can be roughly divided into two main components: data hub (or just “hub”) component and web component. The hub component allows developers to automate the process of monitoring, parsing and uploading your data source to an [Elasticsearch](#) backend. From here, the web component, built on the high-concurrency [Tornado Web Server](#), allows you to easily setup a live high-performance API. The API endpoints expose simple-to-use yet powerful query features using [Elasticsearch's full-text query capabilities and query language](#).

1.3 BioThings API

We also use “*BioThings API*” (or *BioThings APIs*) to refer to an API (or a collection of APIs) built with BioThings SDK. For example, both our popular [MyGene.Info](#) and [MyVariant.Info](#) APIs are built and maintained using this BioThings SDK.

Installing BioThings SDK

You can install the latest stable BioThings SDK release with pip from PyPI, like:

```
pip install biothings
```

You can install the latest development version of BioThings SDK directly from our github repository like:

```
pip install git+https://github.com/biothings/biothings.api.git#egg=biothings
```

Alternatively, you can download the source code, or clone the [BioThings SDK repository](#) and run:

```
python setup.py install
```

2.1 Single Data Source, No Source Updating Tutorial

The following tutorial shows a minimal use case for the BioThings SDK: creating a high-performance, high-concurrency API from a single flat-file. The BioThings SDK is broadly divided into two components, the hub and the web. The hub component is a collection of tools to automate the downloading of source data files, the merging of different sources, and the updating of the Elasticsearch index. The web component is a Tornado-based API app that subsequently serves data from this Elasticsearch index.

2.1.1 Prerequisites

Before starting, there are a few requirements that need to be installed and configured.

Python

The BioThings SDK requires [Python version 3.4 or higher](#) for full functionality. We recommend installing all python dependencies into a [virtualenv](#).

BioThings SDK

Either install from source, like:

```
git clone https://github.com/biothings/biothings.api.git
cd biothings.api
python setup.py install
```

or use pip, like:

```
pip install biothings
```

or directly from our repository, like:

```
pip install git+https://github.com/biothings/biothings.api.git#egg=biothings
```

Elasticsearch

BioThings APIs currently serve data from an Elasticsearch index, so Elasticsearch is a requirement. Install Elasticsearch 2.4 either [directly](#), or as a [docker container](#).

Configure Elasticsearch

To configure Elasticsearch, execute the following commands as su:

```
echo 'http.enabled: True' >> /etc/elasticsearch/elasticsearch.yml
echo 'network.host: "0.0.0.0"' >> /etc/elasticsearch/elasticsearch.yml
```

Note: This guide was created using Ubuntu 16.04, the exact location of elasticsearch.yml may vary in other platforms.

2.1.2 PharmGKB Gene

Once all prerequisites have been installed and Elasticsearch is running, the data loading step can begin. Consider the following script, which defines a “load_data” function that parses the [PharmGKB gene flat file](#) and then iterates through it, storing the results in an Elasticsearch index using `biothings.utils.es.ESIndexer`.

```
In [1]: import re

In [2]: from biothings.utils.es import ESIndexer

In [3]: def load_data(f):
...:     for (i, line) in enumerate(f):
...:         line = line.strip('\n')
...:         if i == 0: # get the column header names in the first row
...:             header_dict = dict([(p, re.sub(r'\s', '_', h.lower())) for (p, h) in
↪in enumerate(line.split('\t'))])
...:         else:
...:             _r = {}
...:             for (pos, val) in enumerate(line.split('\t')):
...:                 if val:
...:                     _r[header_dict[pos]] = val if "','" not in val else val.
↪strip('\"').split('\"', '\"')
```



```

...:         yield _r
...:
In [4]: indexer = ESIndexer(index='pharmgkb_gene_current', doc_type='pharmgkb_gene',
↳ es_host='localhost:9200')

In [5]: indexer.create_index(mapping={'pharmgkb_gene':{'dynamic': True}})

In [6]: with open('genes.tsv', 'r') as gene_file:
...:     indexer.index_bulk(load_data(gene_file))

```

2.1.3 Generate BioThings API

Now that we have an Elasticsearch index with our indexed gene data in it, we can create and start an API. Change to a directory you want to store the front-end code, and type:

```
biothings-admin.py pharmgkb_gene . -o src_package=pharmgkb_gene
```

Now you can start your API by typing:

```
cd pharmgkb_gene/src
pip install -r ../requirements_web.txt
python web/index.py --debug --port=8001
```

Your API is live. To use it, you can query it with a curl (or your local browser). For example, if you wanted to find the PharmGKB accession for an NCBI gene (or gene list) you have, you could do a query like:

```
curl "http://localhost:8001/v1/query?q=ncbi_gene_id:1017&fields=pharmgkb_accession_id"
{
  "max_score": 8.178926,
  "took": 9,
  "total": 1,
  "hits": [
    {
      "_id": "AVydiHIJYMgArMwkfE8R",
      "_score": 8.178926,
      "pharmgkb_accession_id": "PA101"
    }
  ]
}
```

Or, to find all PharmGKB genes that have a CDK* symbol, you can do this query:

```
curl "http://localhost:8001/v1/query?q=symbol:CDK*&fields=pharmgkb_accession_id,symbol"
↳ "
{
  "max_score": 1.0,
  "took": 11,
  "total": 50,
  "hits": [
    {
      "_id": "AVydiHIJYMgArMwkfE8F",
      "_score": 1.0,
      "pharmgkb_accession_id": "PA99",
      "symbol": "CDK1"
    },

```

```
{
  {
    "_id": "AVydiHIJYMgArMwkfE8H",
    "_score": 1.0,
    "pharmgkb_accession_id": "PA26263",
    "symbol": "CDK11A"
  },
  {
    "_id": "AVydiHIJYMgArMwkfE8M",
    "_score": 1.0,
    "pharmgkb_accession_id": "PA165696414",
    "symbol": "CDK15"
  },
  {
    "_id": "AVydiHIJYMgArMwkfE8R",
    "_score": 1.0,
    "pharmgkb_accession_id": "PA101",
    "symbol": "CDK2"
  },
  {
    "_id": "AVydiHIJYMgArMwkfE8n",
    "_score": 1.0,
    "pharmgkb_accession_id": "PA26317",
    "symbol": "CDKL1"
  },
  {
    "_id": "AVydiHIJYMgArMwkfE8N",
    "_score": 1.0,
    "pharmgkb_accession_id": "PA33095",
    "symbol": "CDK16"
  },
  {
    "_id": "AVydiHIJYMgArMwkfE8e",
    "_score": 1.0,
    "pharmgkb_accession_id": "PA38632",
    "symbol": "CDK5RAP2"
  },
  {
    "_id": "AVydiHIJYMgArMwkfE8h",
    "_score": 1.0,
    "pharmgkb_accession_id": "PA26314",
    "symbol": "CDK7"
  },
  {
    "_id": "AVydiHIJYMgArMwkfE8m",
    "_score": 1.0,
    "pharmgkb_accession_id": "PA134871999",
    "symbol": "CDKAL1"
  },
  {
    "_id": "AVydiHIJYMgArMwkfE8v",
    "_score": 1.0,
    "pharmgkb_accession_id": "PA106",
    "symbol": "CDKN2A"
  }
}
]
```

2.2 Multiple Data Sources, Automated Source Updating Tutorial

In this tutorial, we will build the whole process, or “hub”, which produces the data for Taxonomy BioThings API, accessible at t.biothings.io. This API serves information about species, lineage, etc... This “hub” is used to download, maintain up-to-date, process, merge data. At the end of this process, an Elasticsearch index is created containing all the data of interest, ready to be served as an API, using Biothings SDK Web component (covered in another tutorial). Taxonomy Biothings API code is available at <https://github.com/biothings/biothings.species>

2.2.1 Prerequisites

BioThings SDK uses MongoDB as the “staging” storage backend for JSON objects before they are sent to Elasticsearch for indexing. You must have a working MongoDB instance you can connect to. We’ll also perform some basic commands.

You also have to install the latest stable BioThings SDK release, with pip from PyPI:

```
pip install biothings
```

You can install the latest development version of BioThings SDK directly from our github repository like:

```
pip install git+https://github.com/biothings/biothings.api.git#egg=biothings
```

Alternatively, you can download the source code, or clone the [BioThings SDK repository](#) and run:

```
python setup.py install
```

You may want to use `virtualenv` to isolate your installation.

Finally, BioThings SDK is written in python, so you must know some basics.

2.2.2 Configuration file

Before starting to implement our hub, we first need to define a configuration file. A `config_common.py` [file](https://github.com/biothings/biothings.species/blob/master/src/config_common.py) contains all the required configuration variables, some **have** to be defined in your own application, other **can** be overridden as needed.

Typically we will have to define the following:

- MongoDB connections parameters, `DATA_SRC_*` and `DATA_TARGET_*` parameters. They define connections to two different databases, one will contain individual collections for each datasource (SRC) and the other will contain merged collections (TARGET).
- `HUB_DB_BACKEND` defines a database connection for hub purpose (application specific data, like sources status, etc...). Default backend type is MongoDB. We will need to provide a valid `mongodb://` URI. Other backend types are available, like `sqlite3` and `ElasticSearch`, but since we’ll use MongoDB to store and process our data, we’ll stick to the default.
- `DATA_ARCHIVE_ROOT` contains the path of the root folder that will contain all the downloaded and processed data. Other parameters should be self-explanatory and probably don’t need to be changed.
- `LOG_FOLDER` contains the log files produced by the hub

Create a `config.py` and add `from config_common import *` then define all required variables above. `config.py` will look something like this:

```
from config_common import *

DATA_SRC_SERVER = "myhost"
DATA_SRC_PORT = 27017
DATA_SRC_DATABASE = "tutorial_src"
DATA_SRC_SERVER_USERNAME = None
DATA_SRC_SERVER_PASSWORD = None

DATA_TARGET_SERVER = "myhost"
DATA_TARGET_PORT = 27017
DATA_TARGET_DATABASE = "tutorial"
DATA_TARGET_SERVER_USERNAME = None
DATA_TARGET_SERVER_PASSWORD = None

HUB_DB_BACKEND = {
    "module" : "biothings.utils.mongo",
    "uri" : "mongodb://myhost:27017",
}

DATA_ARCHIVE_ROOT = "/tmp/tutorial"
LOG_FOLDER = "/tmp/tutorial/logs"
```

Note: Log folder must be created manually

2.2.3 hub.py

This script represents the main hub executable. Each hub should define it, this is where the different hub commands are going to be defined and where tasks are actually running. It's also from this script that a SSH server will run so we can actually log into the hub and access those registered commands.

Along this tutorial, we will enrich that script. For now, we're just going to define a JobManager, the SSH server and make sure everything is running fine.

```
import asyncio, asyncssh, sys
import concurrent.futures
from functools import partial

import config, biothings
biothings.config_for_app(config)

from biothings.utils.manager import JobManager

loop = asyncio.get_event_loop()
process_queue = concurrent.futures.ProcessPoolExecutor(max_workers=2)
thread_queue = concurrent.futures.ThreadPoolExecutor()
loop.set_default_executor(process_queue)
jmanager = JobManager(loop,
                      process_queue, thread_queue,
                      max_memory_usage=None,
                      )
```

jmanager is our JobManager, it's going to be used everywhere in the hub, each time a parallelized job is created. Species hub is a small one, there's no need for many process workers, two should be fine.

Next, let's define some basic commands for our new hub:

```

from biotthings.utils.hub import schedule, top, pending, done
COMMANDS = {
    "sch" : partial(schedule,loop),
    "top" : partial(top,process_queue,thread_queue),
    "pending" : pending,
    "done" : done,
}

```

These commands are then registered in the SSH server, which is linked to a python interpreter. Commands will be part of the interpreter's namespace and be available from a SSH connection.

```

passwords = {
    'guest': '', # guest account with no password
}

from biotthings.utils.hub import start_server
server = start_server(loop, "Taxonomy hub",passwords=passwords,port=7022,
↳commands=COMMANDS)

try:
    loop.run_until_complete(server)
except (OSError, asyncssh.Error) as exc:
    sys.exit('Error starting server: ' + str(exc))

loop.run_forever()

```

Let's try to run that script ! The first run, it will complain about some missing SSH key:

```

AssertionError: Missing key 'bin/ssh_host_key' (use: 'ssh-keygen -f bin/ssh_host_key'
↳to generate it

```

Let's generate it, following instruction. Now we can run it again and try to connect:

```

$ ssh guest@localhost -p 7022
The authenticity of host '[localhost]:7022 ([127.0.0.1]:7022)' can't be established.
RSA key fingerprint is SHA256:USgdr9nlFVryr475+kQWlLyPxwzIUREcnOCyctUlylQ.
Are you sure you want to continue connecting (yes/no)? yes
Warning: Permanently added '[localhost]:7022' (RSA) to the list of known hosts.

Welcome to Taxonomy hub, guest!
hub>

```

Let's try a command:

```

hub> top()
0 running job(s)
0 pending job(s), type 'top(pending)' for more

```

Nothing fancy here, we don't have much in our hub yet, but everything is running fine.

2.2.4 Dumpers

BioThings species API gathers data from different datasources. We will need to define different dumpers to make this data available locally for further processing.

Taxonomy dumper

This dumper will download taxonomy data from NCBI FTP server. There's one file to download, available at this location: <ftp://ftp.ncbi.nih.gov/pub/taxonomy/taxdump.tar.gz>.

When defining a dumper, we'll need to choose a base class to derive our dumper class from. There are different base dumper classes available in BioThings SDK, depending on the protocol we want to use to download data. In this case, we'll derive our class from `biothings.dataload.dumper.FTPDumper`. In addition to defining some specific class attributes, we will need to implement a method called `create_todump_list()`. This method fills `self.to_dump` list, which is later going to be used to download data. One element in that list is a dictionary with the following structure:

```
{"remote": "<path to file on remote server", "local": "<local path to file>"}
```

Remote information are relative to the working directory specified as class attribute. Local information is an absolute path, containing filename used to save data.

Let's start coding. We'll save that python module in `dataload/sources/taxonomy/dumper.py`.

```
import biothings, config
biothings.config_for_app(config)
```

Those lines are used to configure BioThings SDK according to our own configuration information.

```
from config import DATA_ARCHIVE_ROOT
from biothings.dataload.dumper import FTPDumper
```

We then import a configuration constant, and the FTPDumper base class.

```
class TaxonomyDumper(FTPDumper):

    SRC_NAME = "taxonomy"
    SRC_ROOT_FOLDER = os.path.join(DATA_ARCHIVE_ROOT, SRC_NAME)
    FTP_HOST = 'ftp.ncbi.nih.gov'
    CWD_DIR = '/pub/taxonomy'
    SUFFIX_ATTR = "timestamp"
    SCHEDULE = "0 9 * * *"
```

- `SRC_NAME` will be used as the registered name for this datasource (more on this later).
- `SRC_ROOT_FOLDER` is the folder path for this resource, without any version information (dumper will create different sub-folders for each version).
- `FTP_HOST` and `CWD_DIR` gives information to connect to the remote FTP server and move to appropriate remote directory (`FTP_USER` and `FTP_PASSWD` constants can also be used for authentication).
- `SUFFIX_ATTR` defines the attributes that's going to be used to create folder for each downloaded version. It's basically either "release" or "timestamp", depending on whether the resource we're trying to dump has an actual version. Here, for taxdump file, there's no version, so we're going to use "timestamp". This attribute is automatically set to current date, so folders will look like that: `.../taxonomy/20170120`, `.../taxonomy/20170121`, etc...
- Finally `SCHEDULE`, if defined, will allow that dumper to regularly run within the hub. This is a cron-like notation (see `aiocron` documentation for more).

We now need to tell the dumper what to download, that is, create that `self.to_dump` list:

```
def create_todump_list(self, force=False):
    file_to_dump = "taxdump.tar.gz"
```

```

new_localfile = os.path.join(self.new_data_folder, file_to_dump)
try:
    current_localfile = os.path.join(self.current_data_folder, file_to_dump)
except TypeError:
    # current data folder doesn't even exist
    current_localfile = new_localfile
if force or not os.path.exists(current_localfile) or self.remote_is_better(file_
↳to_dump, current_localfile):
    # register new release (will be stored in backend)
    self.to_dump.append({"remote": file_to_dump, "local":new_localfile})

```

That method tries to get the latest downloaded file and then compare that file with the remote file using `self.remote_is_better(file_to_dump, current_localfile)`, which compares the dates and return True if the remote is more recent. A dict is then created with required elements and appended to `self.to_dump` list.

When the dump is running, each element from that `self.to_dump` list will be submitted to a job and be downloaded in parallel. Let's try our new dumper. We need to update `hub.py` script to add a `DumperManager` and then register this dumper:

In `hub.py`:

```

import dataload
import biothings.dataload.dumper as dumper

dmanager = dumper.DumperManager(job_manager=jmanager)
dmanager.register_sources(dataload.__sources__)
dmanager.schedule_all()

```

Let's also register new commands in the hub:

```

COMMANDS = {
    # dump commands
    "dm" : dmanager,
    "dump" : dmanager.dump_src,
    ...
}

```

`dm` will a shortcut for the dumper manager object, and `dump` will actually call manager's `dump_src()` method.

Manager is auto-registering dumpers from list defines in `dataload` package. Let's define that list:

In `dataload/__init__.py`:

```

__sources__ = [
    "dataload.sources.taxonomy",
]

```

That's it, it's just a string pointing to our taxonomy package. We'll expose our dumper class in that package so the manager can inspect it and find our dumper (note: we could use give the full path to our dumper module, `dataload.sources.taxonomy.dumper`, but we'll add uploaders later, it's better to have one single line per resource).

In `dataload/sources/taxonomy/__init__.py`

```

from .dumper import TaxonomyDumper

```

Let's run the hub again. We can on the logs that our dumper has been found:

```

Found a class based on BaseDumper: '<class 'dataload.sources.taxonomy.dumper.
↳TaxonomyDumper'>'

```

Also, manager has found scheduling information and created a task for this:

```
Scheduling task functools.partial(<bound method DumperManager.create_and_dump of
↳<DumperManager [1 registered]: ['taxonomy']>>, <class 'dataload.sources.taxonomy.
↳dumper.TaxonomyDumper'>, job_manager=<biothings.utils.manager.JobManager object at 0x7f88fc5346d8>, force=False): 0 9 * * *
```

We can double-check this by connecting to the hub, and type some commands:

```
Welcome to Taxonomy hub, guest!
hub> dm
<DumperManager [1 registered]: ['taxonomy']>
```

When printing the manager, we can check our taxonomy resource has been registered properly.

```
hub> sch()
DumperManager.create_and_dump(<class 'dataload.sources.taxonomy.dumper.TaxonomyDumper
↳'>,) [0 9 * * * ] {run in 00h:39m:09s}
```

Dumper is going to run in 39 minutes ! We can trigger a manual upload too:

```
hub> dump("taxonomy")
[1] RUN {0.0s} dump("taxonomy")
```

OK, dumper is running, we can follow task status from the console. At some point, task will be done:

```
hub>
[1] OK dump("taxonomy"): finished, [None]
```

It successfully run (OK), nothing was returned by the task ([None]). Logs show some more details:

```
DEBUG:taxonomy.hub:Creating new TaxonomyDumper instance
INFO:taxonomy_dump:1 file(s) to download
DEBUG:taxonomy_dump:Downloading 'taxdump.tar.gz'
INFO:taxonomy_dump:taxonomy successfully downloaded
INFO:taxonomy_dump:success
```

Alright, now if we try to run the dumper again, nothing should be downloaded since we got the latest file available. Let's try that, here are the logs:

```
DEBUG:taxonomy.hub:Creating new TaxonomyDumper instance
DEBUG:taxonomy_dump:'taxdump.tar.gz' is up-to-date, no need to download
INFO:taxonomy_dump:Nothing to dump
```

So far so good! The actual file, depending on the configuration settings, it's located in `./data/taxonomy/20170125/taxdump.tar.gz`. We can notice the timestamp used to create the folder. Let's also have a look at in the internal database to see the resource status. Connect to MongoDB:

```
> use hub_config
switched to db hub_config
> db.src_dump.find()
{
  "_id" : "taxonomy",
  "release" : "20170125",
  "data_folder" : "./data/taxonomy/20170125",
  "pending_to_upload" : true,
  "download" : {
    "logfile" : "./data/taxonomy/taxonomy_20170125_dump.log",
```



```

        "time" : "4.52s",
        "status" : "success",
        "started_at" : ISODate("2017-01-25T08:32:28.448Z")
    }
}
>

```

We have some information about the download process, how long it took to download files, etc... We have the path to the `data_folder` containing the latest version, the `release` number (here, it's a timestamp), and a flag named `pending_to_upload`. That will be used later to automatically trigger an upload after a dumper has run.

So the actual file is currently compressed, we need to uncompress it before going further. We can add a post-dump step to our dumper. There are two options there, by overriding one of those methods:

```

def post_download(self, remotefile, localfile): triggered for each downloaded file
def post_dump(self): triggered once all files have been downloaded

```

We could use either, but there's a utility function available in BioThings SDK that uncompress everything in a directory, let's use it in a global post-dump step:

```

from biothings.utils.common import untargzall
...

def post_dump(self):
    untargzall(self.new_data_folder)

```

`self.new_data_folder` is the path to the folder freshly created by the dumper (in our case, `./data/taxonomy/20170125`)

Let's try this in the console (restart the hub to make those changes alive). Because file is up-to-date, dumper will not run. We need to force it:

```

hub> dump("taxonomy", force=True)

```

Or, instead of downloading the file again, we can directly trigger the post-dump step:

```

hub> dump("taxonomy", steps="post")

```

There are 2 steps available in a dumper:

1. **dump** : will actually download files
2. **post** : will post-process downloaded files (`post_dump`)

By default, both run sequentially.

After typing either of these commands, logs will show some information about the uncompressing step:

```

DEBUG:taxonomy.hub:Creating new TaxonomyDumper instance
INFO:taxonomy_dump:success
INFO:root:untargz '/opt/slelong/Documents/Projects/biothings.species/src/data/
↳taxonomy/20170125/taxdump.tar.gz'

```

Folder contains all uncompressed files, ready to be process by an uploader.

UniProt species dumper

Following guideline from previous taxonomy dumper, we're now implementing a new dumper used to download species list. There's just one file to be downloaded from ftp://ftp.uniprot.org/pub/databases/uniprot/current_release/

knowledgebase/complete/docs/specplist.txt. Same as before, dumper will inherit FTPDumper base class. File is not compressed, so except this, this dumper will look the same.

Code is available on github for further details: [ee674c55bad849b43c8514fcc6b7139423c70074](#) for the whole commit changes, and [dataload/sources/uniprot/dumper.py](#) for the actual dumper.

Gene information dumper

The last dumper we have to implement will download some gene information from NCBI (ftp://ftp.ncbi.nlm.nih.gov/gene/DATA/gene_info.gz). It's very similar to the first one (we could even have merged them together).

Code is available on github: [d3b3486f71e865235efd673d2f371b53eaa0bc5b](#) for whole changes and [dataload/sources/geneinfo/dumper.py](#) for the dumper.

2.2.5 Uploaders

Now that we have local data available, we can process them. We're going to create 3 different uploaders, one for each datasource. Each uploader will load data into MongoDB, into individual/single collections. Those will then be used in the last merging step.

Before going further, we'll first create an UploaderManager instance and register some of its commands in the hub:

```
import biotthings.dataload.uploader as uploader
# will check every 10 seconds for sources to upload
umanager = uploader.UploaderManager(poll_schedule = '* * * * */10', job_
↳manager=jmanager)
umanager.register_sources(dataload.__sources__)
umanager.poll()

COMMANDS = {
...
    # upload commands
    "um" : umanager,
    "upload" : umanager.upload_src,
...
}
```

Running the hub, we'll see the kind of log statements:

```
INFO:taxonomy.hub:Found 2 resources to upload (['species', 'geneinfo'])
INFO:taxonomy.hub:Launch upload for 'species'
ERROR:taxonomy.hub:Resource 'species' needs upload but is not registered in manager
INFO:taxonomy.hub:Launch upload for 'geneinfo'
ERROR:taxonomy.hub:Resource 'geneinfo' needs upload but is not registered in manager
...
```

Indeed, datasources have been dumped, and a `pending_to_upload` flag has been to `True` in `src_dump`. Upload-Manager polls this `src_dump` internal collection, looking for this flag. If set, it runs automatically the corresponding uploader(s). Since we didn't implement any uploaders yet, manager complains... Let's fix that.

Taxonomy uploader

The taxonomy files we downloaded need to be parsed and stored into a MongoDB collection. We won't go in too much details regarding the actual parsing, there are two parsers, one for `nodes.dmp` and another for `names.dmp` files. They yield dictionaries as the result of this parsing step. We just need to "connect" those parsers to uploaders.

Following the same approach as for dumpers, we’re going to implement our first uploaders by inheriting one the base classes available in BioThings SDK. We have two files to parse, data will stored in two different MongoDB collections, so we’re going to have two uploaders. Each inherits from `biothings.dataload.uploader.BaseSourceUploader`, `load_data` method has to be implemented, this is where we “connect” parsers.

Beside this method, another important point relates to the storage engine. `load_data` will, through the parser, yield documents (dictionaries). This data is processed internally by the base uploader class (`BaseSourceUploader`) using a storage engine. `BaseSourceUploader` uses `biothings.dataload.storage.BasicStorage` as its engine. This storage inserts data in MongoDB collection using bulk operations for better performances. There are other storages available, depending on how data should be inserted (eg. `IgnoreDuplicatedStorage` will ignore any duplicated data error). While choosing a base uploader class, we need to consider which storage class it’s actually using behind-the-scene (an alternative way to do this is using `BaseSourceUploader` and set the class attribute `storage_class`, such as in this uploader: [biothings/dataload/uploader.py#L447](#)).

The first uploader will take care of `nodes.dmp` parsing and storage.

```
import biothings.dataload.uploader as uploader
from .parser import parse_refseq_names, parse_refseq_nodes

class TaxonomyNodesUploader(uploader.BaseSourceUploader):

    main_source = "taxonomy"
    name = "nodes"

    def load_data(self, data_folder):
        nodes_file = os.path.join(data_folder, "nodes.dmp")
        self.logger.info("Load data from file '%s'" % nodes_file)
        return parse_refseq_nodes(open(nodes_file))
```

- `TaxonomyNodesUploader` derives from `BaseSourceUploader`
- `name` gives the name of the collection used to store the data. If `main_source` is *not* defined, it must match `SRC_NAME` in dumper’s attributes
- `main_source` is optional and allows to define main sources and sub-sources. Since we have 2 parsers here, we’re going to have 2 collections created. For this one, we want the collection named “nodes”. But this parser relates to *taxonomy* datasource, so we define a main source called **taxonomy**, which matches `SRC_NAME` in dumper’s attributes.
- `load_data()` has `data_folder` as parameter. It will be set accordingly, to the path of the last version dumped. Also, that method gets data from parsing function `parse_refseq_nodes`. It’s where we “connect” the parser. We just need to return parser’s result so the storage can actually store the data.

The other parser, for `names.dmp`, is almost the same:

```
class TaxonomyNamesUploader(uploader.BaseSourceUploader):

    main_source = "taxonomy"
    name = "names"

    def load_data(self, data_folder):
        names_file = os.path.join(data_folder, "names.dmp")
        self.logger.info("Load data from file '%s'" % names_file)
        return parse_refseq_names(open(names_file))
```

We then need to “expose” those parsers in `taxonomy` package, in `dataload/sources/taxonomy/__init__.py`:

```
from .uploader import TaxonomyNodesUploader, TaxonomyNamesUploader
```

Now let's try to run the hub again. We should see uploader manager has automatically triggered some uploads:

```
INFO:taxonomy.hub:Launch upload for 'taxonomy'
...
...
INFO:taxonomy.names_upload:Uploading 'names' (collection: names)
INFO:taxonomy.nodes_upload:Uploading 'nodes' (collection: nodes)
INFO:taxonomy.nodes_upload:Load data from file './data/taxonomy/20170125/nodes.dmp'
INFO:taxonomy.names_upload:Load data from file './data/taxonomy/20170125/names.dmp'
INFO:root:Uploading to the DB...
INFO:root:Uploading to the DB...
```

While running, we can check what jobs are running, using top() command:

```
hub> top()
  PID | SOURCE | CATEGORY | STEP |
  --+--|-----+-----+-----+-----+
  ↳ DESCRIPTION | MEM | CPU | STARTED_AT | DURATION |
5795 | taxonomy.nodes | uploader | update_data |
  ↳ | 49.7MiB | 0.0% | 2017/01/25 14:58:40 | 15.49s |
5796 | taxonomy.names | uploader | update_data |
  ↳ | 54.6MiB | 0.0% | 2017/01/25 14:58:40 | 15.49s |
2 running job(s)
0 pending job(s), type 'top(pending)' for more
16 finished job(s), type 'top(done)' for more
```

We can see two uploaders running at the same time, one for each file. top (done) can also display jobs that are done and finally top (pending) can give an overview of jobs that are going to be launched when a worker is available (it happens when there are more jobs created than the available number of workers overtime).

In src_dump collection, we can see some more information about the resource and its upload processes. Two jobs were created, we have information about the duration, log files, etc...

```
> db.src_dump.find({_id:"taxonomy"})
{
  "_id" : "taxonomy",
  "download" : {
    "started_at" : ISODate("2017-01-25T13:09:26.423Z"),
    "status" : "success",
    "time" : "3.31s",
    "logfile" : "./data/taxonomy/taxonomy_20170125_dump.log"
  },
  "data_folder" : "./data/taxonomy/20170125",
  "release" : "20170125",
  "upload" : {
    "status" : "success",
    "jobs" : {
      "names" : {
        "started_at" : ISODate("2017-01-25T14:58:40.034Z"),
        "pid" : 5784,
        "logfile" : "./data/taxonomy/taxonomy.names_20170125_
↳upload.log",
        "step" : "names",
        "temp_collection" : "names_temp_eJUdh1te",
        "status" : "success",
        "time" : "26.61s",
        "count" : 1552809,
        "time_in_s" : 27
      },
      "nodes" : {
```

```

        "started_at" : ISODate("2017-01-25T14:58:40.043Z"),
        "pid" : 5784,
        "logfile" : "./data/taxonomy/taxonomy.nodes_20170125_
↪upload.log",
        "step" : "nodes",
        "temp_collection" : "nodes_temp_T5VnzRQC",
        "status" : "success",
        "time" : "22.4s",
        "time_in_s" : 22,
        "count" : 1552809
    }
}
}
}

```

In the end, two collections were created, containing parsed data:

```

> db.names.count()
1552809
> db.nodes.count()
1552809

> db.names.find().limit(2)
{
  "_id" : "1",
  "taxid" : 1,
  "other_names" : [
    "all"
  ],
  "scientific_name" : "root"
}
{
  "_id" : "2",
  "other_names" : [
    "bacteria",
    "not bacteria haeckel 1894"
  ],
  "genbank_common_name" : "eubacteria",
  "in-part" : [
    "monera",
    "procaryotae",
    "prokaryota",
    "prokaryotae",
    "prokaryote",
    "prokaryotes"
  ],
  "taxid" : 2,
  "scientific_name" : "bacteria"
}

> db.nodes.find().limit(2)
{ "_id" : "1", "rank" : "no rank", "parent_taxid" : 1, "taxid" : 1 }
{
  "_id" : "2",
  "rank" : "superkingdom",
  "parent_taxid" : 131567,
  "taxid" : 2
}

```

UniProt species uploader

Following the same guideline, we're going to create another uploader for species file.

```
import biotthings.dataupload.uploader as uploader
from .parser import parse_uniprot_speclist

class UniprotSpeciesUploader(uploader.BaseSourceUploader):

    name = "uniprot_species"

    def load_data(self, data_folder):
        nodes_file = os.path.join(data_folder, "speclist.txt")
        self.logger.info("Load data from file '%s'" % nodes_file)
        return parse_uniprot_speclist(open(nodes_file))
```

In that case, we need only one uploader, so we just define “name” (no need to define `main_source` here).

We need to expose that uploader from the package, in `dataupload/sources/uniprot/__init__.py`:

```
from .uploader import UniprotSpeciesUploader
```

Let's run this through the hub. We can use the “upload” command there (though manager should trigger the upload itself):

```
hub> upload("uniprot_species")
[1] RUN {0.0s} upload("uniprot_species")
```

Similar to dumpers, there are different steps we can individually call for an uploader:

- **data**: will take care of storing data
- **post**: calls `post_update()` method, once data has been inserted. Useful to post-process data or create an index for instance
- **master**: will register the source in `src_master` collection, which is used during the merge step. Uploader method `get_mapping()` can optionally returns an ElasticSearch mapping, it will be stored in `src_master` during that step. We'll see more about this later.
- **clean**: will clean temporary collections and other leftovers...

Within the hub, we can specify these steps manually (they're all executed by default).

```
hub> upload("uniprot_species", steps="clean")
```

Or using a list:

```
hub> upload("uniprot_species", steps=["data", "clean"])
```

Gene information uploader

Let's move forward and implement the last uploader. The goal for this uploader is to identify whether, for a taxonomy ID, there are existing/known genes. File contains information about genes, first column is the `taxid`. We want to know all taxonomy IDs present in the file, and the merged document, we want to add key such as `{ 'has_gene' : True/False }`.

Obviously, we're going to have a lot of duplicates, because for one `taxid` we can have many genes present in the files. We have options here 1) remove duplicates before inserting data in database, or 2) let the database handle the

duplicates (rejecting them). Though we could process data in memory – processed data is rather small in the end –, for demo purpose, we'll go for the second option.

```
import biorthings.dataload.uploader as uploader
import biorthings.dataload.storage as storage
from .parser import parse_geneinfo_taxid

class GeneInfoUploader(uploader.BaseSourceUploader):

    storage_class = storage.IgnoreDuplicatedStorage

    name = "geneinfo"

    def load_data(self, data_folder):
        gene_file = os.path.join(data_folder, "gene_info")
        self.logger.info("Load data from file '%s'" % gene_file)
        return parse_geneinfo_taxid(open(gene_file))
```

- `storage_class`: this is the most important setting in this case, we want to use a storage that will ignore any duplicated records.
- `parse_geneinfo_taxid`: is the parsing function, yield documents as `{"_id" : "taxid"}`

The rest is closed to what we already encountered. Code is available on github in [dataload/sources/geneinfo/uploader.py](https://github.com/biothings/dataload/sources/geneinfo/uploader.py)

When running the uploader, logs show statements like these:

```
INFO:taxonomy.hub:Found 1 resources to upload (['geneinfo'])
INFO:taxonomy.hub:Launch upload for 'geneinfo'
INFO:taxonomy.hub:Building task: functools.partial(<bound method UploaderManager.
↳ create_and_load of <UploaderManager [3 registered]: ['geneinfo', 'species',
↳ 'taxonomy']>>, <class 'dataload.sources.gen
info.uploader.GeneInfoUploader'>, job_manager=<biothings.utils.manager.JobManager_
↳ object at 0x7fbf5f8c69b0>)
INFO:geneinfo_upload:Uploading 'geneinfo' (collection: geneinfo)
INFO:geneinfo_upload:Load data from file './data/geneinfo/20170125/gene_info'
INFO:root:Uploading to the DB...
INFO:root:Inserted 62 records, ignoring 9938 [0.3s]
INFO:root:Inserted 15 records, ignoring 9985 [0.28s]
INFO:root:Inserted 0 records, ignoring 10000 [0.23s]
INFO:root:Inserted 31 records, ignoring 9969 [0.25s]
INFO:root:Inserted 16 records, ignoring 9984 [0.26s]
INFO:root:Inserted 4 records, ignoring 9996 [0.21s]
INFO:root:Inserted 4 records, ignoring 9996 [0.25s]
INFO:root:Inserted 1 records, ignoring 9999 [0.25s]
INFO:root:Inserted 26 records, ignoring 9974 [0.23s]
INFO:root:Inserted 61 records, ignoring 9939 [0.26s]
INFO:root:Inserted 77 records, ignoring 9923 [0.24s]
```

While processing data in batch, some are inserted, others (duplicates) are ignored and discarded. The file is quite big, so the process can be long...

Note: should we want to implement the first option, the parsing function would build a dictionary indexed by taxid and would read the whole, extracting taxid. The whole dict would then be returned, and then processed by storage engine.

So far, we've defined dumpers and uploaders, made them working together through some managers defined in the hub. We're now ready to move the last step: merging data.

2.2.6 Mergers

Merging will be the last step in our hub definition. So far we have data about species, taxonomy and whether a taxonomy ID has known genes in NCBI. In the end, we want to have a collection where documents look like this:

```
{
  _id: "9606",
  authority: ["homo sapiens linnaeus, 1758"],
  common_name: "man",
  genbank_common_name: "human",
  has_gene: true,
  lineage: [9606, 9605, 207598, 9604, 314295, 9526, ...],
  other_names: ["humans"],
  parent_taxid: 9605,
  rank: "species",
  scientific_name: "homo sapiens",
  taxid: 9606,
  uniprot_name: "homo sapiens"
}
```

- `_id`: the taxid, the ID used in all of our individual collection, so the key will be used to collect documents and merge them together (it's actually a requirement, documents are merged using `_id` as the common key).
- `authority`, `common_name`, `genbank_common_name`, `other_names`, `scientific_name` and `taxid` come from `taxonomy.names` collection.
- `uniprot_name` comes from `species` collection.
- `has_gene` is a flag set to true, because taxid 9606 has been found in collection `geneinfo`.
- `parent_taxid` and `rank` come from `taxonomy.nodes` collection.
- (there can be other fields available, but basically the idea here is to merge all our individual collections...)
- finally, `lineage`... it's a little tricky as we need to query nodes in order to compute that field from `_id` and `parent_taxid`.

A first step would be to merge **names**, **nodes** and **species** collections together. Other keys need some post-merge processing, they will be handled in a second part.

Let's first define a `BuilderManager` in the hub.

```
import biotthings.databuild.builder as builder
bmanager = builder.BuilderManager(poll_schedule='* * * * */10', job_
↳manager=jmanager)
bmanager.configure()
bmanager.poll()

COMMANDS = {
...
  # building/merging
  "bm" : bmanager,
  "merge" : bmanager.merge,
...
}
```

Merging configuration

`BuilderManager` uses a builder class for merging. While there are many different dumpers and uploaders classes, there's only one merge class (for now). The merging process is defined in a configuration collection named `src_build`.

Usually, we have as many configurations as merged collections, in our case, we'll just define one configuration.

When running the hub with a builder manager registered, manager will automatically create this `src_build` collection and create configuration placeholder.

```
> db.src_build.find()
{
  "_id" : "placeholder",
  "name" : "placeholder",
  "sources" : [ ],
  "root" : [ ]
}
```

We're going to use that template to create our own configuration:

- **_id** and **name** are the name of the configuration (they can be different but really, **_id** is the one used here)... We'll set these as: `{"_id": "mytaxonomy", "name": "mytaxonomy" }`.
- **sources** is a list of collection names used for the merge. A element is this can also be a regular expression matching collection names. If we have data spread across different collection, like one collection per chromosome data, we could use a regex such as `data_chr.*`. We'll set this as: `{"sources": ["names", "species", "nodes", "geneinfo"]}`
- **root** defines root datasources, that is, datasources that can be used to initiate document creation. Sometimes, we want data to be merged only if an existing document previously exists in the merged collection. If root sources are defined, they will be merged first, then the other remaining in sources will be merged with existing documents. If root doesn't exist (or list is empty), all sources can initiate documents creation. **root** can be a list of collection names, or a negation (not a mix of both). So, for instance, if we want all datasources to be root, except `source10`, we can specify: `"root" : ["!source10"]`. Finally, all root sources must all be declared in sources (root is a subset of sources). That said, it's interesting in our case because we have taxonomy information coming from NCBI and UniProt, but we want to make sure a document built from UniProt only doesn't exist (it's because we need `parent_taxid` field which only exists in NCBI data, so we give priority to those sources first). So root sources are going to be `names` and `nodes`, but because we're lazy typist, we're going to set this to: `{"root" : ["!species"]}`

The resulting document should look like this. Let's save this in `src_build` (and also remove the placeholder, not useful anymore):

```
> conf
{
  "_id" : "mytaxonomy",
  "name" : "mytaxonomy",
  "sources" : [
    "names",
    "uniprot_species",
    "nodes",
    "geneinfo"
  ],
  "root" : ["!uniprot_species"]
}
> db.src_build.save(conf)
> db.src_build.remove({"_id": "placeholder"})
```

Note: **geneinfo** contains only IDs, we could ignore it while merging but we'll need it to be declared as a source when we'll create the index later.

Restarting the hub, we can then check that configuration has properly been registered in the manager, ready to be used. We can list the sources specified in configuration.

```
hub> bm
<BuilderManager [1 registered]: ['mytaxonomy']>
hub> bm.list_sources("mytaxonomy")
['names', 'species', 'nodes']
```

OK, let's try to merge !

```
hub> merge("mytaxonomy")
[1] RUN {0.0s} merge("mytaxonomy")
```

Looking at the logs, we can see builder will first root sources:

```
INFO:mytaxonomy_build:Merging into target collection 'mytaxonomy_20170127_pnlygtqp'
...
INFO:mytaxonomy_build:Sources to be merged: ['names', 'nodes', 'species', 'geneinfo']
INFO:mytaxonomy_build:Root sources: ['names', 'nodes', 'geneinfo']
INFO:mytaxonomy_build:Other sources: ['species']
INFO:mytaxonomy_build:Merging root document sources: ['names', 'nodes', 'geneinfo']
```

Then once root sources are processed, **species** collection merged on top on existing documents:

```
INFO:mytaxonomy_build:Merging other resources: ['species']
DEBUG:mytaxonomy_build:Documents from source 'species' will be stored only if a
↳previous document exists with same _id
```

After a while, task is done, merge has returned information about the amount of data that have been merge: 1552809 records from collections **names**, **nodes** and **geneinfo**, 25394 from **species**. Note: the figures show the number fetched from collections, not necessarily the data merged. For instance, merged data from **species** may be less since it's not a root datasource).

```
hub>
[1] OK merge("mytaxonomy"): finished, [{'total_species': 25394, 'total_nodes': 1552809, 'total_names': 1552809}]
```

Builder creates multiple merger jobs per collection. The merged collection name is, by default, generating from the build name (**mytaxonomy**), and contains also a timestamp and some random chars. We can specify the merged collection name from the hub. By default, all sources defined in the configuration are merged., and we can also select one or more specific sources to merge:

```
hub> merge("mytaxonomy", sources="uniprot_species", target_name="test_merge")
```

Note: sources parameter can also be a list of string.

If we go back to `src_build`, we can have information about the different merges (or builds) we ran:

```
> db.src_build.find({_id:"mytaxonomy"},{build:1})
{
  "_id" : "mytaxonomy",
  "build" : [
    ...
    {
      "src_versions" : {
        "geneinfo" : "20170125",
        "taxonomy" : "20170125",
        "uniprot_species" : "20170125"
      },
      "time_in_s" : 386,
      "logfile" : "./data/logs/mytaxonomy_20170127_build.log",
```

```

    "pid" : 57702,
    "target_backend" : "mongo",
    "time" : "6m26.29s",
    "step_started_at" : ISODate("2017-01-27T11:36:47.401Z"),
    "stats" : {
      "total_uniprot_species" : 25394,
      "total_nodes" : 1552809,
      "total_names" : 1552809
    },
    "started_at" : ISODate("2017-01-27T11:30:21.114Z"),
    "status" : "success",
    "target_name" : "mytaxonomy_20170127_pnlygtqp",
    "step" : "post-merge",
    "sources" : [
      "uniprot_species"
    ]
  }
}

```

We can see the merged collection (auto-generated) is **mytaxonomy_20170127_pnlygtqp**. Let's have a look at the content (remember, collection is in target database, not in src):

```

> use tutorial
switched to db tutorial
> db.mytaxonomy_20170127_pnlygtqp.count()
1552809
> db.mytaxonomy_20170127_pnlygtqp.find({_id:9606})
{
  "_id" : 9606,
  "rank" : "species",
  "parent_taxid" : 9605,
  "taxid" : 9606,
  "common_name" : "man",
  "other_names" : [
    "humans"
  ],
  "scientific_name" : "homo sapiens",
  "authority" : [
    "homo sapiens linnaeus, 1758"
  ],
  "genbank_common_name" : "human",
  "uniprot_name" : "homo sapiens"
}

```

Both collections have properly been merged. We now have to deal with the other data.

Mappers

The next bit of data we need to merge is **geneinfo**. As a reminder, this collection only contains taxonomy ID (as `_id` key) which have known NCBI genes. We'll create a mapper, containing this information. A mapper basically acts as an object that can pre-process documents while they are merged.

Let's define that mapper in `databuild/mapper.py`

```

import biorthings, config
biorthings.config_for_app(config)
from biorthings.utils.common import loadobj
import biorthings.utils.mongo as mongo

```

```

import biotthings.databuild.mapper as mapper
# just to get the collection name
from dataload.sources.geneinfo.uploader import GeneInfoUploader

class HasGeneMapper(mapper.BaseMapper):

    def __init__(self, *args, **kwargs):
        super(HasGeneMapper, self).__init__(*args, **kwargs)
        self.cache = None

    def load(self):
        if self.cache is None:
            # this is a whole dict containing all taxonomy_ids
            col = mongo.get_src_db()[GeneInfoUploader.name]
            self.cache = [d["_id"] for d in col.find({}, {"_id":1})]

    def process(self, docs):
        for doc in docs:
            if doc["_id"] in self.cache:
                doc["has_gene"] = True
            else:
                doc["has_gene"] = False
            yield doc

```

We derive our mapper from `biotthings.databuild.mapper.BaseMapper`, which expects `load` and `process` methods to be defined. `load` is automatically called when the mapper is used by the builder, and `process` contains the main logic, iterating over documents, optionally enrich them (it can also be used to filter documents, by not yielding them). The implementation is pretty straightforward. We get and cache the data from `geneinfo` collection (the whole collection is very small, less than 20'000 IDs, so it can fit nicely and efficiently in memory). If a document has its `_id` found in the cache, we enrich it.

Once defined, we register that mapper into the builder. In `bin/hub.py`, we modify the way we define the builder manager:

```

import biotthings.databuild.builder as builder
from databuild.mapper import HasGeneMapper
hasgene = HasGeneMapper(name="has_gene")
pbuilder = partial(builder.DataBuilder, mappers=[hasgene])
bmanager = builder.BuilderManager(
    poll_schedule='* * * * * */10',
    job_manager=jmanager,
    builder_class=pbuilder)
bmanager.configure()
bmanager.poll()

```

First we instantiate a mapper object and give it a name (more on this later). While creating the manager, we need to pass a builder class. The problem here is we also have to give our mapper to that class while it's instantiated. We're using `partial` (from `functools`), which allows to partially define the class instantiation. In the end, `builder_class` parameter is expected to a callable, which is the case with `partial`.

Let's try if our mapper works (restart the hub). Inside the hub, we're going to manually get a builder instance. Remember through the SSH connection, we can access python interpreter's namespace, which is very handy when it comes to develop and debug as we can directly access and "play" with objects and their states:

First we get a builder instance from the manager:

```
hub> builder = bm["mytaxonomy"]
hub> builder
<biothings.databuild.builder.DataBuilder object at 0x7f278aecf400>
```

Let's check the mappers and get ours:

```
hub> builder.mappers
{None: <biothings.databuild.mapper.TransparentMapper object at 0x7f278aecf4e0>, 'has_
↪gene': <databuild.mapper.HasGeneMapper object at 0x7f27ac6c0a90>}
```

We have our `has_gene` mapper (it's the name we gave). We also have a `TransparentMapper`. This mapper is automatically added and is used as the default mapper for any document (there has to be one...).

```
hub> hasgene = builder.mappers["has_gene"]
hub> len(hasgene.cache)
Error: TypeError("object of type 'NoneType' has no len()",)
```

Oops, cache isn't loaded yet, we have to do it manually here (but it's done automatically during normal execution).

```
hub> hasgene.load()
hub> len(hasgene.cache)
19201
```

OK, it's ready. Let's now talk more about the mapper's name. A mapper can be applied to different sources, and we have to define which sources' data should go through that mapper. In our case, we want **names** and **species** collection's data to go through. In order to do that, we have to instruct the uploader with a special attribute. Let's modify `dataload.sources.species.uploader.UniprotSpeciesUploader` class

```
class UniprotSpeciesUploader(uploader.BaseSourceUploader):

    name = "uniprot_species"
    __metadata__ = {"mapper" : 'has_gene'}
```

`__metadata__` dictionary is going to be used to create a master document. That document is stored in `src_master` collection (we talked about it earlier). Let's add this metadata to `dataload.sources.taxonomy.uploader.TaxonomyNamesUploader`

```
class TaxonomyNamesUploader(uploader.BaseSourceUploader):

    main_source = "taxonomy"
    name = "names"
    __metadata__ = {"mapper" : 'has_gene'}
```

Before using the builder, we need to refresh master documents so these metadata are stored in `src_master`. We could trigger a new upload, or directly tell the hub to only process master steps:

```
hub> upload("uniprot_species", steps="master")
[1] RUN {0.0s} upload("uniprot_species", steps="master")
hub> upload("taxonomy.names", steps="master")
[1] OK upload("uniprot_species", steps="master"): finished, [None]
[2] RUN {0.0s} upload("taxonomy.names", steps="master")
```

(you'll notice for taxonomy, we only trigger upload for sub-source **names**, using "dot-notation", corresponding to "main_source.name". Let's now have a look at those master documents:

```
> db.src_master.find({'_id':{'$in':['uniprot_species', 'names']}})
{
```

```

    "_id" : "names",
    "name" : "names",
    "timestamp" : ISODate("2017-01-26T16:21:32.546Z"),
    "mapper" : "has_gene",
    "mapping" : {

    }
}
{
    "_id" : "uniprot_species",
    "name" : "uniprot_species",
    "timestamp" : ISODate("2017-01-26T16:21:19.414Z"),
    "mapper" : "has_gene",
    "mapping" : {

    }
}

```

We have our mapper key stored. We can now trigger a new merge (we specify the target collection name):

```

hub> merge("mytaxonomy",target_name="mytaxonomy_test")
[3] RUN {0.0s} merge("mytaxonomy",target_name="mytaxonomy_test")

```

In the logs, we can see our mapper has been detected and is used:

```

INFO:mytaxonomy_build:Creating merger job #1/16, to process 'names' 100000/1552809 (6.
↪4%)
INFO:mytaxonomy_build:Found mapper '<databuild.mapper.HasGeneMapper object at_
↪0x7f47ef3bbac8>' for source 'names'
INFO:mytaxonomy_build:Creating merger job #1/1, to process 'species' 25394/25394 (100.
↪0%)
INFO:mytaxonomy_build:Found mapper '<databuild.mapper.HasGeneMapper object at_
↪0x7f47ef3bbac8>' for source 'species'

```

Once done, we can query the merged collection to check the data:

```

> use tutorial
switched to db tutorial
> db.mytaxonomy_test.find({_id:9606})
{
  "_id" : "9606",
  "has_gene" : true,
  "taxid" : 9606,
  "uniprot_name" : "homo sapiens",
  "other_names" : [
    "humans"
  ],
  "scientific_name" : "homo sapiens",
  "authority" : [
    "homo sapiens linnaeus, 1758"
  ],
  "genbank_common_name" : "human",
  "common_name" : "man"
}

```

OK, there's a has_gene flag that's been set. So far so good !

Post-merge process

We need to add lineage and parent taxid information for each of these documents. We'll implement that last part as a post-merge step, iterating over each of them. In order to do so, we need to define our own builder class to override proper methods there. Let's define it in `databuild/builder.py`.

```
import biotthings.databuild.builder as builder
import config

class TaxonomyDataBuilder(builder.DataBuilder):

    def post_merge(self, source_names, batch_size, job_manager):
        pass
```

The method we have to implement in `post_merge`, as seen above. We also need to change `hub.py` to use that builder class:

```
from databuild.builder import TaxonomyDataBuilder
pbuilder = partial(TaxonomyDataBuilder, mappers=[hasgene])
```

For now, we just added a class level in the hierarchy, everything runs the same as before. Let's have a closer look to that post-merge process. For each document, we want to build the lineage. Information is stored in **nodes** collection. For instance, taxid 9606 (homo sapiens) has a parent_taxid 9605 (homo), which has a parent_taxid 207598 (homininae), etc... In the end, the homo sapiens lineage is:

```
9606, 9605, 207598, 9604, 314295, 9526, 314293, 376913, 9443, 314146, 1437010,
9347, 32525, 40674, 32524, 32523, 1338369, 8287, 117571, 117570, 7776, 7742,
89593, 7711, 33511, 33213, 6072, 33208, 33154, 2759, 131567 and 1
```

We could recursively query **nodes** collections until we reach the top the tree, but that would be a lot of queries. We just need `taxid` and `parent_taxid` information to build the lineage, maybe it's possible to build a dictionary that could fit in memory. **nodes** has 1552809 records. A dictionary would use $2 * 1552809 * \text{sizeof}(\text{integer}) + \text{index overhead}$. That's probably few megabytes, let's assume that ok... (note: using `pympler` lib, we can actually know that dictionary size will be closed to 200MB...)

We're going to use another mapper here, but no sources will use it. We'll just instantiate it from `post_merge` method. In `databuild/mapper.py`, let's add another class:

```
from dataload.sources.taxonomy.uploader import TaxonomyNodesUploader
```

```
class LineageMapper(mapper.BaseMapper):

    def __init__(self, *args, **kwargs):
        super(LineageMapper, self).__init__(*args, **kwargs)
        self.cache = None

    def load(self):
        if self.cache is None:
            col = mongo.get_src_db()[TaxonomyNodesUploader.name]
            self.cache = {}
            [self.cache.setdefault(d["_id"], d["parent_taxid"]) for d in col.find({}, {
↪ "parent_taxid": 1})]

    def get_lineage(self, doc):
        if doc['taxid'] == doc['parent_taxid']: #take care of node #1
            # we reached the top of the taxonomy tree
            doc['lineage'] = [doc['taxid']]
            return doc
        # initiate lineage with information we have in the current doc
```

```

lineage = [doc['taxid'], doc['parent_taxid']]
while lineage[-1] != 1:
    parent = self.cache[lineage[-1]]
    lineage.append(parent)
doc['lineage'] = lineage
return doc

def process(self, docs):
    for doc in docs:
        doc = self.get_lineage(doc)
        yield doc

```

Let's use that mapper in `TaxonomyDataBuilder`'s `post_merge` method. The signature is the same as `merge()` method (what's actually called from the hub) but we just need the `batch_size` one: we're going to grab documents from the merged collection in batch, process them and update them in batch as well. It's going to be much faster than dealing one document at a time. To do so, we'll use `doc_feeder` utility function:

```

from biothings.utils.mongo import doc_feeder, get_target_db
from biothings.databuild.builder import DataBuilder
from biothings.dataload.storage import UpsertStorage

from databuild.mapper import LineageMapper
import config
import logging

class TaxonomyDataBuilder(DataBuilder):

    def post_merge(self, source_names, batch_size, job_manager):
        # get the lineage mapper
        mapper = LineageMapper(name="lineage")
        # load cache (it's being loaded automatically
        # as it's not part of an upload process
        mapper.load()

        # create a storage to save docs back to merged collection
        db = get_target_db()
        col_name = self.target_backend.target_collection.name
        storage = UpsertStorage(db, col_name)

        for docs in doc_feeder(self.target_backend.target_collection, step=batch_size,
→ inbatch=True):
            docs = mapper.process(docs)
            storage.process(docs, batch_size)

```

Since we're using the mapper manually, we need to load the cache

- `db` and `col_name` are used to create our storage engine. `Builder` has an attribute called `target_backend` (a `biothings.dataload.backend.TargetDocMongoBackend` object) which can be used to reach the collection we want to work with.
- `doc_feeder` iterates over all the collection, fetching documents in batch. `inbatch=True` tells the function to return data as a list (default is a dict indexed by `_id`).
- those documents are processed by our mapper, setting the lineage information and then are stored using our `UpsertStorage` object.

Note: `post_merge` actually runs within a thread, so any calls here won't block the execution (ie. won't block the `asyncio` event loop execution)

Let's run this on our merged collection. We don't want to merge everything again, so we specify the step we're interested in and the actual merged collection (`target_name`)

```
hub> merge("mytaxonomy",steps="post",target_name="mytaxonomy_test") [1] RUN {0.0s}
merge("mytaxonomy",steps="post",target_name="mytaxonomy_test")
```

After a while, process is done. We can test our updated data:

```
> use tutorial
switched to db tutorial
> db.mytaxonomy_test.find({_id:9606})
{
  "_id" : 9606,
  "taxid" : 9606,
  "common_name" : "man",
  "other_names" : [
    "humans"
  ],
  "uniprot_name" : "homo sapiens",
  "rank" : "species",
  "lineage" : [9606,9605,207598,9604,...,131567,1],
  "genbank_common_name" : "human",
  "scientific_name" : "homo sapiens",
  "has_gene" : true,
  "parent_taxid" : 9605,
  "authority" : [
    "homo sapiens linnaeus, 1758"
  ]
}
```

OK, we have new lineage information (truncated for sanity purpose). Merged collection is ready to be used. It can be used for instance to create and send documents to an Elasticsearch database. This is what's actually occurring when creating a BioThings web-service API. That step will be covered in another tutorial.

2.2.7 Indexers

Coming soon!

Full updated and maintained code for this hub is available here: <https://github.com/biothings/biothings.species>

Also, taxonomy BioThings API can be queried as this URL: <http://t.biothings.io>

2.3 BioThings Standalone instances

This step-by-step guide shows how to use Biothings standalone instances. Standalones instances are based on Docker containers and provide and fully pre-configured, ready-to-use Biothings API that can easily be maintained and kept up-to-date. The idea is, for any user, to be able to run her own APIs locally, fulfilling different needs:

- keep all APIs requests private and local to your own server
- enriching existing and publicly available data found on our APIs with some private data
- run API on your own architecture to perform heavy queries that would sometimes be throttled out from online services

2.3.1 Quick Links

If you already know how to run a BioThings standalone instance, you can download the latest available Docker images from the following tables.

Note: Images don't contain data but are ready to download and maintain data up-to-date running simple commands through the hub.



Production and old data require at least 30GiB disk space.

Production	Demo	Old
Contact us	Download	Download



Production and old data require at least 2TiB disk space.

Production	Demo	Old
Contact us	Download	Download



Production and old data require at least 150Gib disk space.

Production	Demo	Old
Contact us	Download	Soon !

2.3.2 Prerequisites

Using standalone instances requires to have a Docker server up and running, some basic knowledge about commands to run and use containers. Images have been tested on Docker >=17. Using AWS cloud, you can use our public AMI **biothings_demo_docker** (ami-44865e3c) with Docker pre-configured and ready for standalone demo instances deployment. We recommend using instance type with at least 8GiB RAM, such as `t2.large`. AMI comes with an extra 30GiB EBS volume, which should be enough to deploy any demo instances.

Alternately, you can install your own Docker server (on recent Ubuntu systems, `sudo apt-get install docker.io` is usually enough). You may need to point Docker images directory to a specific hard drive to get enough space, using `-g` option:

```
# /mnt/docker points to a hard drive with enough disk space
sudo echo 'DOCKER_OPTS="-g /mnt/docker"' >> /etc/default/docker
# restart to make this change active
sudo service docker restart
```

Demo instances use very little disk space, as only a small subset of data is available. For instance, myvariant demo only requires ~10GiB to run with demo data up-to-date, including the whole Linux system and all other dependencies. Demo instances provide a quick and easy way to setup a running APIs, without having to deal with some advanced system configurations.

For deployment with production or old data, you may need a large amount of disk space. Refer to the [Quick Links](#) section for more information. Bigger instance types will also be required, and even a full cluster architecture deployment. We'll soon provide guidelines and deployment scripts for this purpose.

2.3.3 What you'll learn

Through this guide, you'll learn:

- how to obtain a Docker image to run your favorite API
- how to run that image inside a Docker container and how to access the web API
- how to connect to the *hub*, a service running inside to container used to interact with the API systems
- how to use that hub, using specific commands, in order to perform update and keep data up-to-date

2.3.4 Data found in standalone instances

All BioThings APIs (mygene.info, myvariant.info, ...) provide data release in different flavors:

- **Production data**, the actual data found on live APIs we, the BioThings team at [SuLab](#), are running and keeping up-to-date on a regular basis. Please contact us if you're interested in obtaining this type of data.
- **Demo data**, a small subset of production data, publicly available
- **Old production data**, an at least one year old production dataset (full), publicly available

The following guide applies to demo data only, though the process would be very similar for other types of data flavors.

2.3.5 Downloading and running a standalone instance

Standalone instances are available as Docker images. For the purpose of this guide, we'll setup an instance running mygene API, containing demo data. Links to standalone demo Docker images, can be found in [Quick links](#) at the beginning of this guide. Use one of these links, or use this [direct link](#) to mygene's demo instance, and download the Docker image file, using your favorite browser or `wget`:

```
$ wget http://biothings-containers.s3-website-us-west-2.amazonaws.com/demo_mygene/
↳demo_mygene.docker
```

You must have a running Docker server in order to use that image. Typing `docker ps` should return all running containers, or at least an empty list as in the following example. Depending on the systems and configuration, you may have to add `sudo` in front of this command to access Docker server.

```
$ docker ps
CONTAINER ID        IMAGE               COMMAND             CREATED
↳ STATUS          PORTS              NAMES
```

Once downloaded, the image can be loaded into the server:

```
$ docker image load < demo_mygene.docker
$ docker image list
REPOSITORY                                TAG
↔ IMAGE ID                                CREATED                                SIZE
demo_mygene                                latest
↔ 15d6395e780c                            6 weeks ago                            1.78GB
```

Image is now loaded, size is ~1.78GiB, it contains no data (yet). An docker container can now be instantiated from that image, to create a BioThings standalone instance, ready to be used.

A standalone instance is a pre-configured system containing several parts. BioThings hub is the system used to interact with BioThings backend and perform operations such as downloading data and create/update ElasticSearch indices. Those indices are used by the actual BioThings web API system to serve data to end-users. The hub can be accessed through a standard SSH connection or through REST API calls. In this guide, we'll use the SSH server.

A BioThings instance expose several services on different ports:

- **80**: BioThings web API port
- **7022**: BioThings hub SSH port
- **7080**: BioThings hub REST API port
- **9200**: ElasticSearch port

We will map and expose those ports to the host server using option `-p` so we can access BioThings services without having to enter the container (eg. hub ssh port here will accessible using port 19022).

```
$ docker run --name demo_mygene -p 19080:80 -p 19200:9200 -p 19022:7022 -p 19090:7080 -d demo_mygene
```

Note: Instance will store ElasticSearch data in `/var/lib/elasticsearch/` directory, and downloaded data and logs in `/data/` directory. Those two locations could require extra disk space, if needed Docker option `-v` can be used to mount a directory from the host, inside the container. Please refer to Docker documentnation.

Let's enter the container to check everything is running fine. Services may take a while, up to 1 min, before fully started. If some services are missing, the troubleshooting section may help.

```
$ docker exec -ti demo_mygene /bin/bash

root@a6a6812e2969:/tmp# netstat -tnlp
Active Internet connections (only servers)
Proto Recv-Q Send-Q Local Address           Foreign Address         State       PID/
↔Program name
tcp        0      0 0.0.0.0:7080            0.0.0.0:*               LISTEN      -
tcp        0      0 0.0.0.0:7022            0.0.0.0:*               LISTEN      -
tcp        0      0 0.0.0.0:80              0.0.0.0:*               LISTEN      25/
↔nginx
tcp        0      0 127.0.0.1:8881          0.0.0.0:*               LISTEN      -
tcp        0      0 127.0.0.1:8882          0.0.0.0:*               LISTEN      -
tcp        0      0 127.0.0.1:8883          0.0.0.0:*               LISTEN      -
tcp        0      0 127.0.0.1:8884          0.0.0.0:*               LISTEN      -
tcp        0      0 127.0.0.1:8885          0.0.0.0:*               LISTEN      -
tcp        0      0 127.0.0.1:8886          0.0.0.0:*               LISTEN      -
tcp        0      0 127.0.0.1:8887          0.0.0.0:*               LISTEN      -
tcp        0      0 127.0.0.1:8888          0.0.0.0:*               LISTEN      -
tcp6       0      0 :::7080                 :::*                    LISTEN      -
tcp6       0      0 :::7022                 :::*                    LISTEN      -
```

tcp6	0	0	:::9200	:::*	LISTEN	-
tcp6	0	0	:::9300	:::*	LISTEN	-

We can see the different BioThings services' ports: 7080, 7022 and 7080. All 888x ports correspond to Tornado instances running behind Nginx port 80. They shouldn't be accessed directly. Ports 9200 and 9300 are Elasticsearch standard ports (9200 one can be used to perform queries directly on ES, if needed)

At this point, the standalone instance is up and running. No data has been downloaded yet, let's see how to populate the BioThings API using the hub.

2.3.6 Updating data using Biothings hub

If the standalone instance has been freshly started, there's no data to be queried by the API. If we make a API call, such as fetching metadata, we'll get an error:

```
# from Docker host
$ curl -v http://localhost:19080/metadata
* Trying 127.0.0.1...
* Connected to localhost (127.0.0.1) port 19080 (#0)
> GET /metadata HTTP/1.1
> Host: localhost:19080
> User-Agent: curl/7.47.0
> Accept: */*
>
< HTTP/1.1 500 Internal Server Error
< Date: Tue, 28 Nov 2017 18:19:23 GMT
< Content-Type: text/html; charset=UTF-8
< Content-Length: 93
< Connection: keep-alive
< Server: TornadoServer/4.5.2
<
* Connection #0 to host localhost left intact
```

This 500 error reflects a missing index (ElasticSearch index, the backend used by BioThings web API). We can have a look at existing indices in Elasticsearch:

```
# from Docker host
$ curl http://localhost:19200/_cat/indices
yellow open hubdb 5 1 0 0 795b 795b
```

There's only one index, hubdb, which is an internal index used by the hub. No index containing actual biological data...

BioThings hub is a service running inside the instance, it can be accessed through a SSH connection, or using REST API calls. For the purpose of the guide, we'll use SSH. Let's connect to the hub (type *yes* to accept the key on first connection):

```
# from Docker host
$ ssh guest@localhost -p 19022
The authenticity of host '[localhost]:19022 ([127.0.0.1]:19022)' can't be established.
RSA key fingerprint is SHA256:j63IEgXc3yJqgv0F4wa35aGliH5YQux84xxABew5AS0.
Are you sure you want to continue connecting (yes/no)? yes
Warning: Permanently added '[localhost]:19022' (RSA) to the list of known hosts.

Welcome to Auto-hub, guest!
hub>
```

We're now connected to the hub, inside a python shell where the application is actually running. Let's see what commands are available:

Warning: the hub console, though accessed through SSH, is **not** a Linux shell (such as *bash*), it's a python interpreter shell.

```
hub> help()

Available commands:

    versions
    check
    info
    download
    apply
    step_update
    update
    help

Type: 'help(command)' for more
```

- `versions()` will display all available data build versions we can download to populate the API
- `check()` will return whether a more recent version is available online
- `info()` will display current local API version, and information about the latest available online
- `download()` will download the data compatible with current local version (but without populating the Elasticsearch index)
- `apply()` will use local data previously downloaded to populate the index
- `step_update()` will bring data release to the next one (one step in versions), compatible with current local version
- `update()` will bring data to the latest available online (using a combination of `download` and `apply` calls)

Note: `update()` is the fastest, easiest and preferred way to update the API. `download`, `apply`, `step_update` are available when it's necessary to bring the API data to a specific version (not the latest one), are considered more advanced, and won't be covered in this guide.

Note: Because the hub console is actually a python interpreter, we call the commands using parenthesis, just like functions or methods. We can also pass arguments when necessary, just like standard python (remember: it **is** python...)

Note: After each command is typed, we need to press "enter" to get either its status (still running) or the result

Let's explore some more.

```
hub> info()
[2] RUN {0.0s} info()
hub>
[2] OK info(): finished
>>> Current local version: 'None'
```

```
>>> Release note for remote version 'latest':
Build version: '20171126'
=====
Previous build version: '20171119'
Generated on: 2017-11-26 at 03:11:51

+-----+-----+-----+-----+-----+
↪-----+
| Updated datasource          | prev. release | new release | prev. # of docs | new # of docs |
↪-----+
| entrez.entrez_gene         | 20171118     | 20171125   | 10,003          |                |
↪10,003 |
| entrez.entrez_refseq       | 20171118     | 20171125   | 10,003          |                |
↪10,003 |
| entrez.entrez_unigene      | 20171118     | 20171125   | 10,003          |                |
↪10,003 |
| entrez.entrez_go           | 20171118     | 20171125   | 10,003          |                |
↪10,003 |
| entrez.entrez_genomic_pos  | 20171118     | 20171125   | 10,003          |                |
↪10,003 |
| entrez.entrez_retired      | 20171118     | 20171125   | 10,003          |                |
↪10,003 |
| entrez.entrez_accession    | 20171118     | 20171125   | 10,003          |                |
↪10,003 |
| generif                    | 20171118     | 20171125   | 10,003          |                |
↪10,003 |
| uniprot                    | 20171025     | 20171122   | 10,003          |                |
↪10,003 |
+-----+-----+-----+-----+-----+
↪-----+

Overall, 9,917 documents in this release
0 document(s) added, 0 document(s) deleted, 130 document(s) updated
```

We can see here we don't have any local data release (Current local version: 'None'), whereas the latest online (at that time) is from November 26th 2017. We can also see the release note with the different changes involved in the release (whether it's a new version, or the number of documents that changed).

```
hub> versions()
[1] RUN {0.0s} versions()
hub>
[1] OK versions(): finished
version=20171003          date=2017-10-05T09:47:59.413191 type=full
version=20171009          date=2017-10-09T14:47:10.800140 type=full
version=20171009.20171015 date=2017-10-19T11:44:47.961731 type=incremental
version=20171015.20171022 date=2017-10-25T13:33:16.154788 type=incremental
version=20171022.20171029 date=2017-11-14T10:34:39.445168 type=incremental
version=20171029.20171105 date=2017-11-06T10:55:08.829598 type=incremental
version=20171105.20171112 date=2017-11-14T10:35:04.832871 type=incremental
version=20171112.20171119 date=2017-11-20T07:44:47.399302 type=incremental
version=20171119.20171126 date=2017-11-27T10:38:03.593699 type=incremental
```

Data comes in two distinct types:

- **full**: this is a full data release, corresponding to an ElasticSearch snapshot, containing all the data
- **incremental** : this is a differential/incremental release, produced by computing the differences between two

consecutive versions. The diff data is then used to patch an existing, compatible data release to bring it to the next version.

So, in order to obtain the latest version, the hub will first find a compatible version. Since it's currently empty (no data), it will use the first **full** release from 20171009, and then apply **incremental** updates sequentially (20171009 . 20171015, then 20171015 . 20171022, then 20171022 . 20171029, etc... up to 20171119 . 20171126).

Let's update the API:

```
hub> update()
[3] RUN {0.0s} update()
hub>
[3] RUN {1.3s} update()
hub>
[3] RUN {2.07s} update()
```

After a while, the API is up-to-date, we can run command `info()` again (it also can be used to track update progress):

```
hub> info()
[4] RUN {0.0s} info()
hub>
[4] OK info(): finished
>>> Current local version: '20171126'
>>> Release note for remote version 'latest':
Build version: '20171126'
=====
Previous build version: '20171119'
Generated on: 2017-11-26 at 03:11:51

+-----+-----+-----+-----+-----+
↔-----+
| Updated datasources | prev. release | new release | prev. # of docs | new # of docs |
+-----+-----+-----+-----+-----+
↔-----+
| entrez.entrez_gene | 20171118 | 20171125 | 10,003 | 10,003 |
↔10,003 |
| entrez.entrez_refseq | 20171118 | 20171125 | 10,003 | 10,003 |
↔10,003 |
| entrez.entrez_unigene | 20171118 | 20171125 | 10,003 | 10,003 |
↔10,003 |
| entrez.entrez_go | 20171118 | 20171125 | 10,003 | 10,003 |
↔10,003 |
| entrez.entrez_genomic_pos | 20171118 | 20171125 | 10,003 | 10,003 |
↔10,003 |
| entrez.entrez_retired | 20171118 | 20171125 | 10,003 | 10,003 |
↔10,003 |
| entrez.entrez_accession | 20171118 | 20171125 | 10,003 | 10,003 |
↔10,003 |
| generif | 20171118 | 20171125 | 10,003 | 10,003 |
↔10,003 |
| uniprot | 20171025 | 20171122 | 10,003 | 10,003 |
↔10,003 |
+-----+-----+-----+-----+-----+
↔-----+

Overall, 9,917 documents in this release
0 document(s) added, 0 document(s) deleted, 130 document(s) updated
```


Local version is 20171126, remote is 20171126, we're up-to-date. We can also use `check()`:

```
hub> check()
[5] RUN {0.0s} check()
hub>
[5] OK  check(): finished
Nothing to dump
```

Nothing to dump means there's no available remote version that can be downloaded. It would otherwise return a version number, meaning we would be able to update the API again using command `update()`.

Press Control-D to exit from the hub console.

Querying ElasticSearch, we can see a new index, named `biothings_current`, has been created and populated:

```
$ curl http://localhost:19200/_cat/indices
green open biothings_current 1 0 14903 0 10.3mb 10.3mb
yellow open hubdb 5 1 2 0 11.8kb 11.8kb
```

We now have a populated API we can query:

```
# from Docker host
# get metadata (note the build_version field)
$ curl http://localhost:19080/metadata
{
  "app_revision": "672d55f2deab4c7c0e9b7249d22ccca58340a884",
  "available_fields": "http://mygene.info/metadata/fields",
  "build_date": "2017-11-26T02:58:49.156184",
  "build_version": "20171126",
  "genome_assembly": {
    "rat": "rn4",
    "nematode": "ce10",
    "fruitfly": "dm3",
    "pig": "susScr2",
    "mouse": "mm10",
    "zebrafish": "zv9",
    "frog": "xenTro3",
    "human": "hg38"
  },
},

# annotation endpoint
$ curl http://localhost:19080/v3/gene/1017?fields=alias,ec
{
  "_id": "1017",
  "_score": 9.268311,
  "alias": [
    "CDKN2",
    "p33 (CDK2)"
  ],
  "ec": "2.7.11.22",
  "name": "cyclin dependent kinase 2"
}

# query endpoint
$ curl http://localhost:19080/v3/query?q=cdk2
{
  "max_score": 310.69254,
  "took": 37,
  "total": 10,
```

```
"hits": [
  {
    "_id": "1017",
    "_score": 310.69254,
    "entrezgene": 1017,
    "name": "cyclin dependent kinase 2",
    "symbol": "CDK2",
    "taxid": 9606
  },
  {
    "_id": "12566",
    "_score": 260.58084,
    "entrezgene": 12566,
    "name": "cyclin-dependent kinase 2",
    "symbol": "Cdk2",
    "taxid": 10090
  },
  ...
]
```

2.3.7 BioThings API with multiple indices

Some APIs use more than one Elasticsearch index to run. For instance, `myvariant.info` uses one index for hg19 assembly, and one index for hg38 assembly. With such APIs, the available commands contain a suffix showing which index (thus, which data release) they relate to. Here's the output of `help()` from `myvariant`'s standalone instance:

```
hub> help()

Available commands:

  versions_hg19
  check_hg19
  info_hg19
  download_hg19
  apply_hg19
  step_update_hg19
  update_hg19
  versions_hg38
  check_hg38
  info_hg38
  download_hg38
  apply_hg38
  step_update_hg38
  update_hg38
  help
```

For instance, `update()` command is now available as `update_hg19()` and `update_hg38()` depending on the assembly.

2.3.8 Troubleshooting

We test and make sure, as much as we can, that standalone images are up-to-date and `hub` is properly running for each data release. But things can still go wrong...

First make sure all services are running. Enter the container and type `netstat -tnlp`, you should see services running on ports (see usual running *services*). If services running on ports 7080 or 7022 aren't running, it means the

hub has not started. If you just started the instance, wait a little more as services may take a while before they're fully started and ready.

If after ~1 min, you still don't see the hub running, log to user `biothings` and check the starting sequence.

Note: Hub is running in a tmux session, under user `biothings`

```
# sudo su - biothings
$ tmux a # recall tmux session

python -m biothings.bin.autohub
(pyenv) biothings@a6a6812e2969:~/mygene.info/src$ python -m biothings.bin.autohub
INFO:root:Hub DB backend: {'module': 'biothings.utils.es', 'host': 'localhost:9200'}
INFO:root:Hub database: hubdb
DEBUG:asyncio:Using selector: EpollSelector
start
```

You should see something looking like this above. If not, you should see the actual error, and depending on the error, you may be able to fix it (not enough disk space, etc...). The hub can be started again using `python -m biothings.bin.autohub` from within the application directory (in our case, `/home/biothings/mygene.info/src/`)

Note: Press Control-B then D to detach the tmux session and let the hub running in background.

Logs are available in `/data/mygene.info/logs/`. You can have a look at:

- `dump_*.log` files for logs about data download
- `upload_*.log` files for logs about index update in general (full/incremental)
- `sync_*.log` files for logs about incremental update only
- and `hub_*.log` files for general logs about the hub process

Finally, you can report issues and request for help, by joining Biothings Google Groups (<https://groups.google.com/forum/#!forum/biothings>)

2.4 Hub component

The purpose of the BioThings hub component is to allow you to easily automate the parsing and uploading of your data to an Elasticsearch backend.

2.4.1 dumper

BaseDumper

FTPDumper

HTTPDumper

GoogleDriveDumper

WgetDumper

DummyDumper

ManualDumper

2.4.2 uploader

BaseSourceUploader

NoBatchIgnoreDuplicatedSourceUploader

IgnoreDuplicatedSourceUploader

MergerSourceUploader

DummySourceUploader

ParallelizedSourceUploader

NoDataSourceUploader

2.4.3 builder

DataBuilder

2.4.4 indexer

Indexer

2.4.5 differ

BaseDiffer

JsonDiffer

SelfContainedJsonDiffer

DiffReportRendererBase

2.5. Web component
DiffReportTxt

2.4.6 syncer

2.5.1 Server boot script

BioThings API ioloop start utilities.

This module contains functions to configure and start the [base event loop](#) from command line args. Command line processing is done using `tornado.options`, with the following arguments defined:

- `port`: the port to start the API on, **default** 8000
- `address`: the address to start the API on, **default** 127.0.0.1
- `debug`: start the API in debug mode, **default** False
- `appdir`: path to API configuration directory, **default**: current working directory

The **main** function is the boot script for all BioThings API webservers.

`index_base.main (APP_LIST, app_settings={}, debug_settings={}, sentry_client_key=None)`
Main ioloop configuration and start

Parameters

- **APP_LIST** – a list of `URLSpec` objects or `(regex, handler_class)` tuples
- **app_settings** – Tornado application settings
- **debug_settings** – Additional application settings for API debug mode
- **sentry_client_key** – Application-specific key for attaching Sentry monitor to the application

2.5.2 Settings

Settings objects used to configure the web API. These settings get passed into the `handler.initialize()` function, of each request, and configure the web API endpoint. They are mostly a container for the [Config module](#), and any other settings that are the same across all handler types, e.g. the Elasticsearch client.

Config module

BiothingWebSettings

class `biothings.web.settings.BiothingWebSettings (config='biothings.web.settings.default')`
A container for the settings that configure the web API

The `config` init parameter specifies a module that configures this biothing. For more information see [config module](#) documentation.

generate_app_list ()

Generates the `tornado.web.Application (regex, handler_class, options)` tuples for this project.

set_debug_level (`debug=False`)

Set if running API in debug mode. Should be called before passing `self` to handler initialization.

validate ()

Validate these settings

BiothingESWebSettings

class `biothings.web.settings.BiothingESWebSettings` (*config='biothings.web.settings.default'*)
BiothingWebSettings subclass with functions specific to an elasticsearch backend

The `config` init parameter specifies a module that configures this biothing. For more information see *config module* documentation.

doc_url (*bid*)

Function to return a url on this biothing API to the biothing object specified by bid.

get_es_client ()

Get the [Elasticsearch client](#) for this app, only called once on invocation of server.

source_metadata ()

Function to cache return of the source metadata stored in `_meta` of index mappings

2.5.3 Handlers

BaseHandler

class `biothings.web.api.helper.BaseHandler` (*application, request, **kwargs*)

Parent class of all biothings handlers, only direct descendant of `tornado.web.RequestHandler`, contains the common functions in the biothings handler universe:

- return *self* as JSON
- set CORS and caching headers
- typify the URL keyword arguments
- optionally send tracking data to google analytics and integrate with sentry monitor

ga_event_object (*data={}*)

Create the data object for google analytics tracking.

get_query_params ()

Extract, typify, and sanitize the parameters from the URL query string.

initialize (*web_settings*)

Tornado handler `initialize()`, Override to add settings for *this* biothing API. Assumes that the `web_settings` kwarg exists in `APP_LIST`

log_exceptions (*exception_msg=""*)

Logs the current exception in tornado logs and in hipchat room if available. This must be called in an exception handler

return_json (*data, encode=True, indent=None, status_code=200, _format='json'*)

Return passed data object as JSON response. If `jsonp` parameter is set in the request, return a valid JSONP response.

Parameters

- **data** – object to return as JSON
- **encode** – if encode is False, assumes input data is already a JSON encoded string.
- **indent** – number of indents per level in JSON string
- **status_code** – HTTP status code for response
- **_format** – output format - currently either “html” or “json”

set_cacheable (*etag=None*)

set proper header to make the response cacheable. set etag if provided.

support_cors (**args, **kwargs*)

Provide server side support for CORS request.

BaseESRequestHandler

class `biothings.web.api.es.handlers.base_handler.BaseESRequestHandler` (*application, request, **kwargs*)

Parent class of all Elasticsearch-based Request handlers, subclass of *BaseHandler*. Contains handling for Elasticsearch-specific query params (like *fields*, *size*, etc)

get_cleaned_options (*kwargs*)

Separate URL keyword arguments into their functional category. Incoming kwargs are separated into one of 4 categories, depending on how the argument controls the pipeline:

- *control_kwargs* - These are arguments that control the pipeline execution (**raw**, **rawquery**, etc)
- *es_kwargs* - These are arguments that get passed directly to the Elasticsearch client during query
- *esqb_kwargs* - These are arguments that go to the Elasticsearch query builder (**fields**, **size**, etc)
- *transform_kwargs* - These are arguments that go to the Elasticsearch result transformer (**jsonld**, **dotfield**, etc)

initialize (*web_settings*)

Tornado request handler initialization. Initializations common to all Elasticsearch-specific request handlers go here.

return_raw_query_json (*query, status_code=200, _format='json'*)

Return valid JSON if *rawquery* option is selected. This is necessary as queries can span multiple lines (POST)

BiothingHandler

class `biothings.web.api.es.handlers.biothing_handler.BiothingHandler` (*application, request, **kwargs*)

Request handlers for requests to the annotation lookup endpoint

get (*bid=None*)

Handle a GET to the annotation lookup endpoint.

initialize (*web_settings*)

Tornado handler *.initialize()* function for all requests to the annotation lookup endpoint. Here, the allowed arguments are set (depending on the request method) for each kwarg category.

post (*ids=None*)

Handle a POST to the annotation lookup endpoint

QueryHandler

class `biothings.web.api.es.handlers.query_handler.QueryHandler` (*application*,
request,
***kwargs*)

Request handlers for requests to the query endpoint

get ()

Handle a GET to the query endpoint.

initialize (*web_settings*)

Tornado handler `.initialize()` function for all requests to the query endpoint. Here, the allowed arguments are set (depending on the request method) for each kwarg category.

post ()

Handle a POST to the query endpoint.

MetadataHandler

class `biothings.web.api.es.handlers.metadata_handler.MetadataHandler` (*application*,
re-
quest,
***kwargs*)

Request handlers for requests to the metadata endpoint.

get ()

Handle a GET to the metadata endpoint. Also handles `/metadata/fields`.

initialize (*web_settings*)

Tornado handler `.initialize()` function for all requests to the metadata endpoint. Here, the allowed arguments are set (depending on the request method) for each kwarg category.

2.5.4 Elasticsearch Query Builder

class `biothings.web.api.es.query_builder.ESQueries` (*es_kwargs*={})

A very simple class to object-ize Elasticsearch Query DSL. This should be replaced by the official [Elasticsearch equivalent](#). Also contains a simple JSON validator after generating all queries (dump to string and re-read)

bool (*query_kwargs*)

Given *query_kwargs*, validate and return a **bool** query.

match (*query_kwargs*)

Given *query_kwargs*, validate and return a **match** query.

match_all (*query_kwargs*)

Given *query_kwargs*, validate and return a **match_all** query.

multi_match (*query_kwargs*)

Given *query_kwargs*, validate and return a **multi_match** query.

query_string (*query_kwargs*)

Given *query_kwargs*, validate and return a **query_string** query.

raw_query (*raw_query*)

Given *query_kwargs*, validate and return a *raw* query (queries that don't fit the same *query_template*).

```
class biotthings.web.api.es.query_builder.ESQueryBuilder (index, doc_type, options, es_options, scroll_options={}, userquery_dir="", regex_list=[], default_scopes=['_id'])
```

Class to return the correct query given the request endpoint, method, and URL params.

Parameters

- **index** – The Elasticsearch index to run the query on
- **doc_type** – The Elasticsearch document type of the query
- **options** – Options from the URL string relevant to query building
- **es_options** – Options for Elasticsearch query stage
- **scroll_options** – Options for scroll requests
- **regex_list** – A list of (regex, scope) tuples for annotation lookup
- **userquery_dir** – The directory containing user queries for this app
- **default_scopes** – A list representing the default Elasticsearch query scope(s) for this query

add_extra_filters (*q*)

Override me to add more filters

add_query_filters (*_query*)

Given a query, add any other filters

annotation_GET_query (*bid*)

Return an annotation lookup GET query for this *bid*.

Parameters *bid* – Biothing ID, used to lookup the annotation

annotation_POST_query (*bids*)

Return an annotation lookup POST query for these *bids*.

Parameters *bids* – Biothing IDs, used to lookup the annotations

get_query_filters ()

Override me to add more query filters

metadata_query ()

Return a metadata query

query_GET_query (*q*)

Return a query endpoint GET query for this query string *q*.

Parameters *q* – query string specifying the query

query_POST_query (*qs, scopes*)

Return query endpoint POST queries for these query strings *qs*.

Parameters

- **qs** – Query strings to query
- **scopes** – Scope of query strings *qs*

scroll (*scroll_id*)

Return the next batch of results from a *scroll* query.

Parameters *scroll_id* – ID of the batch to yield hits from

2.5.5 Elasticsearch Query

class `biothings.web.api.es.query.ESQuery` (*client, options={}*)

This class contains functions to execute the *query* section of all handler pipelines. The inputs to it are an Elasticsearch client (from *BiothingESWebSettings*), and any options from the URL string. Each handler calls a different query function, though they all do essentially the same thing: get the query generated in the ESQueryBuilder stage of the pipeline (*query_kwargs*), and run it using the supplied Elasticsearch client.

annotation_GET_query (*query_kwargs*)

Given *query_kwargs* from ESQueryBuilder, return results of annotation lookup GET query on ES client.

annotation_POST_query (*query_kwargs*)

Given *query_kwargs* from ESQueryBuilder, return results of annotation lookup POST query on ES client.

get_biothing (*query_kwargs*)

Return a biothing using the Elasticsearch client.get function

metadata_query (*query_kwargs*)

Given *query_kwargs* from ESQueryBuilder, return results of metadata query on ES client.

query_GET_query (*query_kwargs, *args, **kwargs*)

Given *query_kwargs* from ESQueryBuilder, return results of query GET on ES client.

query_POST_query (*query_kwargs, *args, **kwargs*)

Given *query_kwargs* from ESQueryBuilder, return results of query POST on ES client.

scroll (*query_kwargs*)

Given *query_kwargs* from ESQueryBuilder, return results of a scroll on ES client - returns next batch of results.

2.5.6 Elasticsearch Result Transformer

class `biothings.web.api.es.transform.ESResultTransformer` (*options, host, doc_url_function=<function ESResultTransformer.<lambda>>, jsonld_context={}, data_sources={}, output_aliases={}, app_dir="", source_metadata={}*)

Class to transform the results of the Elasticsearch query generated prior in the pipeline. This contains the functions to extract the final document from the elasticsearch query result in *Elasticsearch Query*. This also contains the code to flatten a document (if **dotfield** is True), or to add JSON-LD context to the document (if **jsonld** is True).

Parameters

- **options** – Options from the URL string controlling result transformer
- **host** – Host name (extracted from request), used for JSON-LD address generation
- **doc_url_function** – a function that takes one argument (a biothing id) and returns a URL to that biothing
- **jsonld_context** – JSON-LD context for this app (optional)
- **data_sources** – unused currently (optional)

- **output_aliases** – list of output key names to alias, unused currently (optional)
- **app_dir** – Application directory for this app (used for getting app information in /meta-data)
- **source_metadata** – Metadata object containing source information for `_license` keys

clean_annotation_GET_response (*res*)

Transform the results of a GET to the annotation lookup endpoint.

Parameters **res** – Results from *Elasticsearch Query*.

clean_annotation_POST_response (*bid_list, res, single_hit=True*)

Transform the results of a POST to the annotation lookup endpoint.

Parameters

- **bid_list** – List of biothing id inputs
- **res** – Results from *Elasticsearch Query*
- **single_hit** – If `True`, render queries with 1 result as a dictionary, else as a 1-element list containing a dictionary

clean_metadata_response (*res, fields=False*)

Transform the results of a GET to the metadata endpoint.

Parameters **res** – Results from *Elasticsearch Query*.

clean_query_GET_response (*res*)

Transform the results of a GET to the query endpoint.

Parameters **res** – Results from *Elasticsearch Query*.

clean_query_POST_response (*qlist, res, single_hit=True*)

Transform the results of a POST to the query endpoint.

Parameters

- **qlist** – List of query inputs
- **res** – Results from *Elasticsearch Query*
- **single_hit** – If `True`, render queries with 1 result as a dictionary, else as a 1-element list containing a dictionary

clean_scroll_response (*res*)

Transform the results of a GET to the scroll endpoint

Parameters **res** – Results from *Elasticsearch Query*.

b

`biothings`, 41

`biothings.web.index_base`, 42

`biothings.web.settings`, 42

A

add_extra_filters() (biothings.web.api.es.query_builder.ESQueryBuilder method), 46

add_query_filters() (biothings.web.api.es.query_builder.ESQueryBuilder method), 46

annotation_GET_query() (biothings.web.api.es.query.ESQuery method), 47

annotation_GET_query() (biothings.web.api.es.query_builder.ESQueryBuilder method), 46

annotation_POST_query() (biothings.web.api.es.query.ESQuery method), 47

annotation_POST_query() (biothings.web.api.es.query_builder.ESQueryBuilder method), 46

B

BaseESRequestHandler (class in biothings.web.api.es.handlers.base_handler), 44

BaseHandler (class in biothings.web.api.helper), 43

BiothingESWebSettings (class in biothings.web.settings), 43

BiothingHandler (class in biothings.web.api.es.handlers.biothing_handler), 44

biothings (module), 39, 41

biothings.web.index_base (module), 42

biothings.web.settings (module), 42

BiothingWebSettings (class in biothings.web.settings), 42

bool() (biothings.web.api.es.query_builder.ESQueries method), 45

C

clean_annotation_GET_response() (biothings.web.api.es.transform.ESResultTransformer

method), 48

clean_annotation_POST_response() (biothings.web.api.es.transform.ESResultTransformer method), 48

clean_metadata_response() (biothings.web.api.es.transform.ESResultTransformer method), 48

clean_query_GET_response() (biothings.web.api.es.transform.ESResultTransformer method), 48

clean_query_POST_response() (biothings.web.api.es.transform.ESResultTransformer method), 48

clean_scroll_response() (biothings.web.api.es.transform.ESResultTransformer method), 48

D

doc_url() (biothings.web.settings.BiothingESWebSettings method), 43

E

ESQueries (class in biothings.web.api.es.query_builder), 45

ESQuery (class in biothings.web.api.es.query), 47

ESQueryBuilder (class in biothings.web.api.es.query_builder), 45

ESResultTransformer (class in biothings.web.api.es.transform), 47

G

ga_event_object() (biothings.web.api.helper.BaseHandler method), 43

generate_app_list() (biothings.web.settings.BiothingWebSettings method), 42

get() (biothings.web.api.es.handlers.biothing_handler.BiothingHandler method), 44

get() (biothings.web.api.es.handlers.metadata_handler.MetadataHandler method), 45

get() (biothings.web.api.es.handlers.query_handler.QueryHandler method), 45
 get_biothing() (biothings.web.api.es.query.ESQuery method), 47
 get_cleaned_options() (biothings.web.api.es.handlers.base_handler.BaseESRequestHandler method), 44
 get_es_client() (biothings.web.settings.BiothingESWebSettings method), 43
 get_query_filters() (biothings.web.api.es.query_builder.ESQueryBuilder method), 46
 get_query_params() (biothings.web.api.helper.BaseHandler method), 43
I
 initialize() (biothings.web.api.es.handlers.base_handler.BaseESRequestHandler method), 44
 initialize() (biothings.web.api.es.handlers.biothing_handler.BiothingHandler method), 44
 initialize() (biothings.web.api.es.handlers.metadata_handler.MetadataHandler method), 45
 initialize() (biothings.web.api.es.handlers.query_handler.QueryHandler method), 45
 initialize() (biothings.web.api.helper.BaseHandler method), 43
L
 log_exceptions() (biothings.web.api.helper.BaseHandler method), 43
M
 main() (biothings.web.index_base method), 42
 match() (biothings.web.api.es.query_builder.ESQueries method), 45
 match_all() (biothings.web.api.es.query_builder.ESQueries method), 45
 metadata_query() (biothings.web.api.es.query.ESQuery method), 47
 metadata_query() (biothings.web.api.es.query_builder.ESQueryBuilder method), 46
 MetadataHandler (class in biothings.web.api.es.handlers.metadata_handler), 45
 multi_match() (biothings.web.api.es.query_builder.ESQueries method), 45
P
 post() (biothings.web.api.es.handlers.biothing_handler.BiothingHandler method), 44
 post() (biothings.web.api.es.handlers.query_handler.QueryHandler method), 45
Q
 query_GET_query() (biothings.web.api.es.query.ESQuery method), 47
 query_GET_query() (biothings.web.api.es.query_builder.ESQueryBuilder method), 46
 query_POST_query() (biothings.web.api.es.query.ESQuery method), 47
 query_POST_query() (biothings.web.api.es.query_builder.ESQueryBuilder method), 46
 query_string() (biothings.web.api.es.query_builder.ESQueries method), 45
 QueryHandler (class in biothings.web.api.es.handlers.query_handler), 45
R
 raw_query() (biothings.web.api.es.query_builder.ESQueries method), 45
 return_json() (biothings.web.api.helper.BaseHandler method), 43
 return_raw_query_json() (biothings.web.api.es.handlers.base_handler.BaseESRequestHandler method), 44
S
 scroll() (biothings.web.api.es.query.ESQuery method), 47
 scroll() (biothings.web.api.es.query_builder.ESQueryBuilder method), 46
 set_cacheable() (biothings.web.api.helper.BaseHandler method), 43
 set_debug_level() (biothings.web.settings.BiothingWebSettings method), 42
 source_metadata() (biothings.web.settings.BiothingESWebSettings method), 43
 support_cors() (biothings.web.api.helper.BaseHandler method), 44
V
 validate() (biothings.web.settings.BiothingWebSettings method), 42